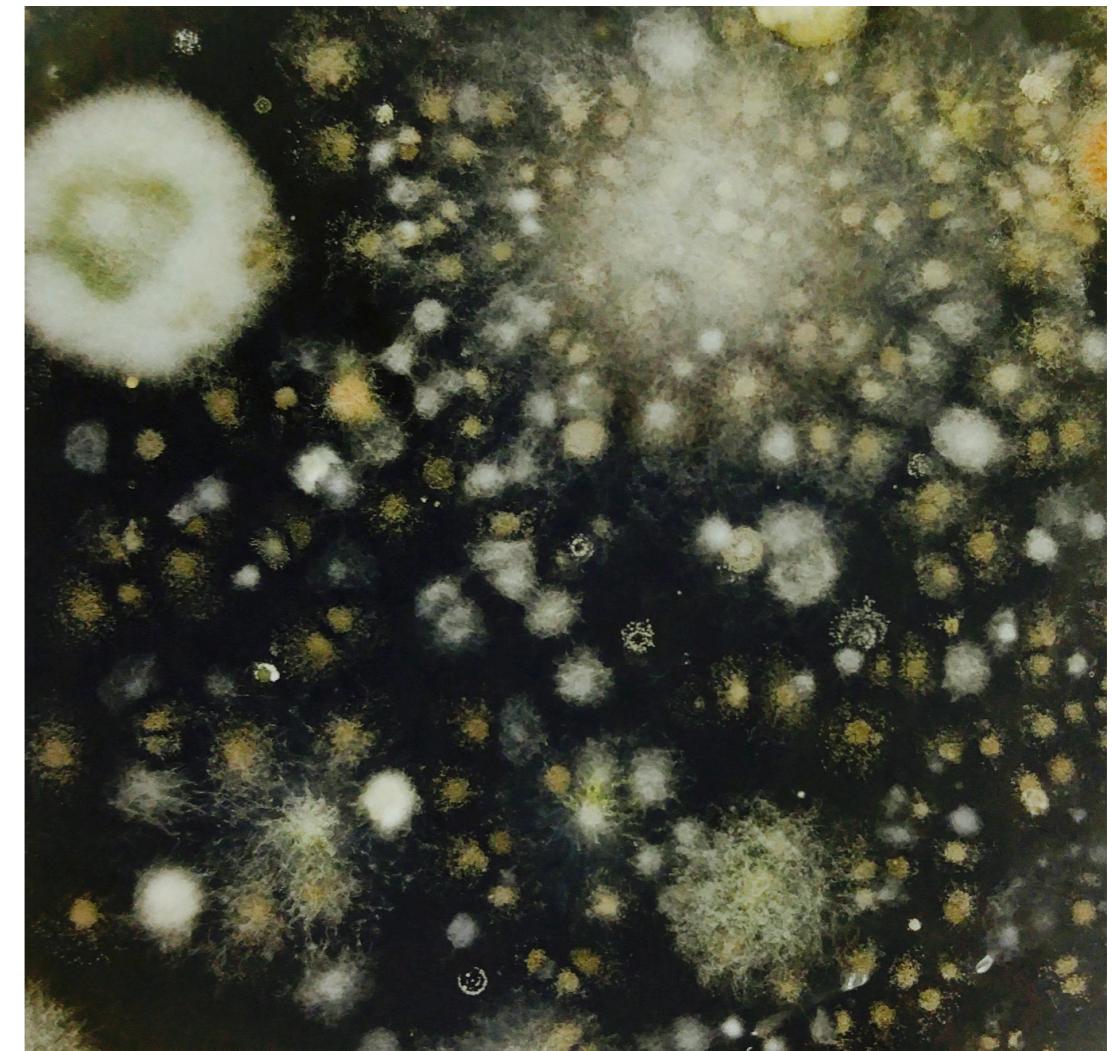


A rigorous & rational approach to

# Microbial differential abundance



Amy D Willis PhD

Associate Professor

Department of Biostatistics

University of Washington

*Pronouns: she/her*

  @AmyDWillis

 [adwillis@uw.edu](mailto:adwillis@uw.edu)



Slides: [github.com/statdivlab/presentations](https://github.com/statdivlab/presentations)

# Paradigms for microbiome data analysis

- Two statistical paradigms:
  1. “How do we *model the data?*”
  2. “How do we *learn about biology?*”

# Paradigm 1: Model the data

- “Modeling the data” conversations sound like
  - compositional? overdispersed? zero-inflated? Negative binomial? Multinomial-Dirichlet? ...
  - rarefy?
  - best  $\alpha$ - and  $\beta$ -diversity metrics?
  - best transformation / “normalization”?

# Paradigm 2: Learn about biology

- “Learning about biology” conversations sound like
  - What exists in the environment?
  - What would we like to learn?
  - How does our data reflect that environment?
  - What can we learn from our data? Under what assumptions?

# Today: Paradigm 2 for differential abundance



Search...  
Help | Adv

arXiv > stat > arXiv:2402.05231

Statistics > Methodology

[Submitted on 7 Feb 2024 (v1), last revised 14 Mar 2025 (this version, v2)]

## Estimating Fold Changes from Partially Observed Outcomes with Applications in Microbial Metagenomics

David S Clausen, Sarah Teichman, Amy D Willis

We consider the problem of estimating fold-changes in the expected value of a multivariate outcome observed with unknown sample-specific and category-specific perturbations. This challenge arises in high-throughput sequencing studies of the abundance of microbial taxa because microbes are systematically over- and under-detected relative to their true abundances. Our model admits a partially identifiable estimand, and we establish full identifiability by imposing interpretable parameter constraints. To reduce bias and guarantee

# What exists in the environment?

- "There is some number of a given biological quantity in every environment"
  - "There are 54,601 *S. epidermidis* cells on my index finger"
  - "There are 0 transcripts of the gene *Core RC1 subunit PsaA* on this podium"
  - "There are 874,455,469 genomes circulating in 100 mL seawater with the 16S variant CGGAGGGTGCA..."

# What exists in the environment?

$Y_{ij}$  = true number of genomic unit  $j$  in environment  $i$

$X_i \in \mathbb{R}^p$  covariate information  
(treatment vs control, age, sex, diet...)

🐱	$Y_{ij}$	💰	I	2	...	J
	ENV I					
	ENV 2					
	...					
	ENV M					
	ENV M+I					
	...					
	ENV N-I					
	ENV N					

# What would we like to know?

$Y_{ij}$  = true number of unit  $j$  in sample  $i$

We do **not** observe  $\{Y_{ij}\} \dots$

...but if we did, what would we do?

🐱 $Y_{ij}$ 💰   2 ... J
ENV I
ENV 2
...
ENV M
ENV M+I
...
ENV N-I
ENV N

# What would we like to know?

- How *abundant* are species?
  - Average of  $Y_{i4}$  across environments
- How *present* are species?
  - % of environments in which  $Y_{i2} > 0$
- How *diverse* are communities?
  - $\#\{j : Y_{ij} > 0\}$
  - $-\sum_{j=1}^J p_{ij} \log p_{ij}$  for  $p_{ij} := \frac{Y_{ij}}{\sum_j Y_{ij}}$
  - ...

🐱 $Y_{ij}$ 💰   1 2 ... J
ENV I
ENV 2
...
ENV M
ENV M+I
...
ENV N-I
ENV N

# What would we like to know?

- Which species differ in their average abundance?
  - All?
- By how much?
  - 0.1%? 50%? 500%?

🐱	$Y_{ij}$	💰	I	2	...	J
ENV I						
ENV 2						
...						
ENV M						
ENV M+I						
...						
ENV N-I						
ENV N						

# What data do we have?

$Y_{ij}$  = true number of unit  $j$  in sample  $i$

$W_{ij}$  = number of times unit  $j$  observed in sample  $i$  from HTS

rainy day icon	$W_{ij}$	crying cat icon		2	...	J
SAMPLE I						
SAMPLE 2						
...						
SAMPLE M						
SAMPLE M+1						
...						
SAMPLE N-I						
SAMPLE N						

How do we connect the  $Y_{ij}$ 's and the  $W_{ij}$ 's?

# Connecting data to reality

- Traditionally, DA methods assume

$$\mathbb{E}[W_{ij}] = c_i Y_{ij}$$

- Is this reasonable?

# Connecting data to reality

- Mock community: An artificially constructed community of known composition

	L.crispatus	L.iners	A.vaginae	S.agalactiae	G.vaginalis	S.amnii	P.bivia
1	0	0	0	1.00	0	0	0
2	0	0	0.5	0	0	0	0.5
3	0.33	0.33	0	0	0	0	0.33
4	0.33	0.33	0	0.33	0	0	0

# Connecting data to reality

	L.crispatus	L.iners	A.vaginae	S.agalactiae	G.vaginalis	S.amnii	P.bivia
1	19	4	2	51332	1	14	1
2	0	1	1424	0	0	7	21708
3	4775	11234	0	0	0	1	3249
4	1644	5497	1	4521	0	7	0

	L.crispatus	L.iners	A.vaginae	S.agalactiae	G.vaginalis	S.amnii	P.bivia
1	0	0	0	1.00	0	0	0
2	0	0	0.5	0	0	0	0.5
3	0.33	0.33	0	0	0	0	0.33
4	0.33	0.33	0	0.33	0	0	0

# Connecting data to reality

1. Despite equal mixing fractions, some taxa are observed many more times
2. Despite being purportedly absent, taxa are observed

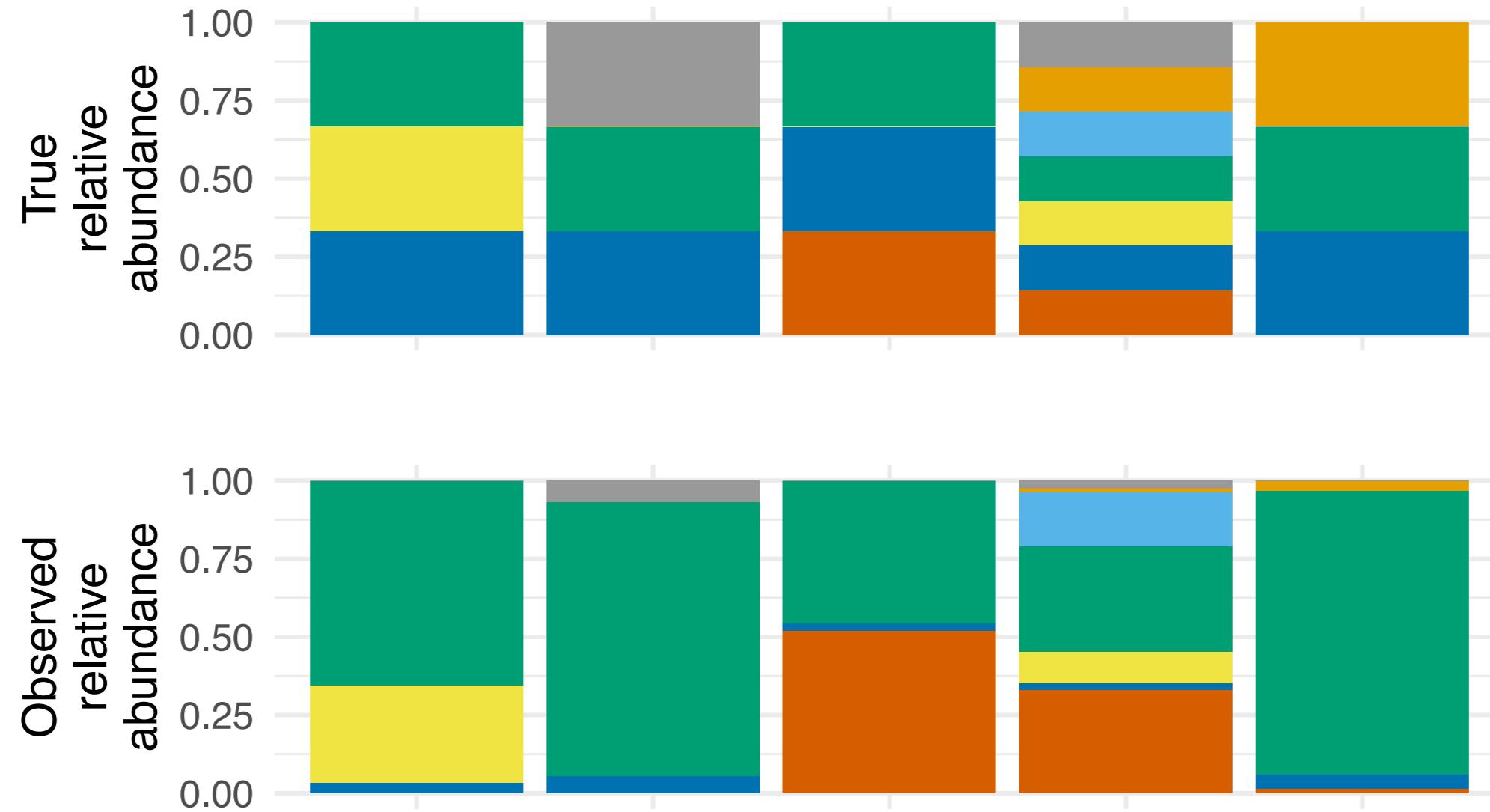
# Connecting data to reality

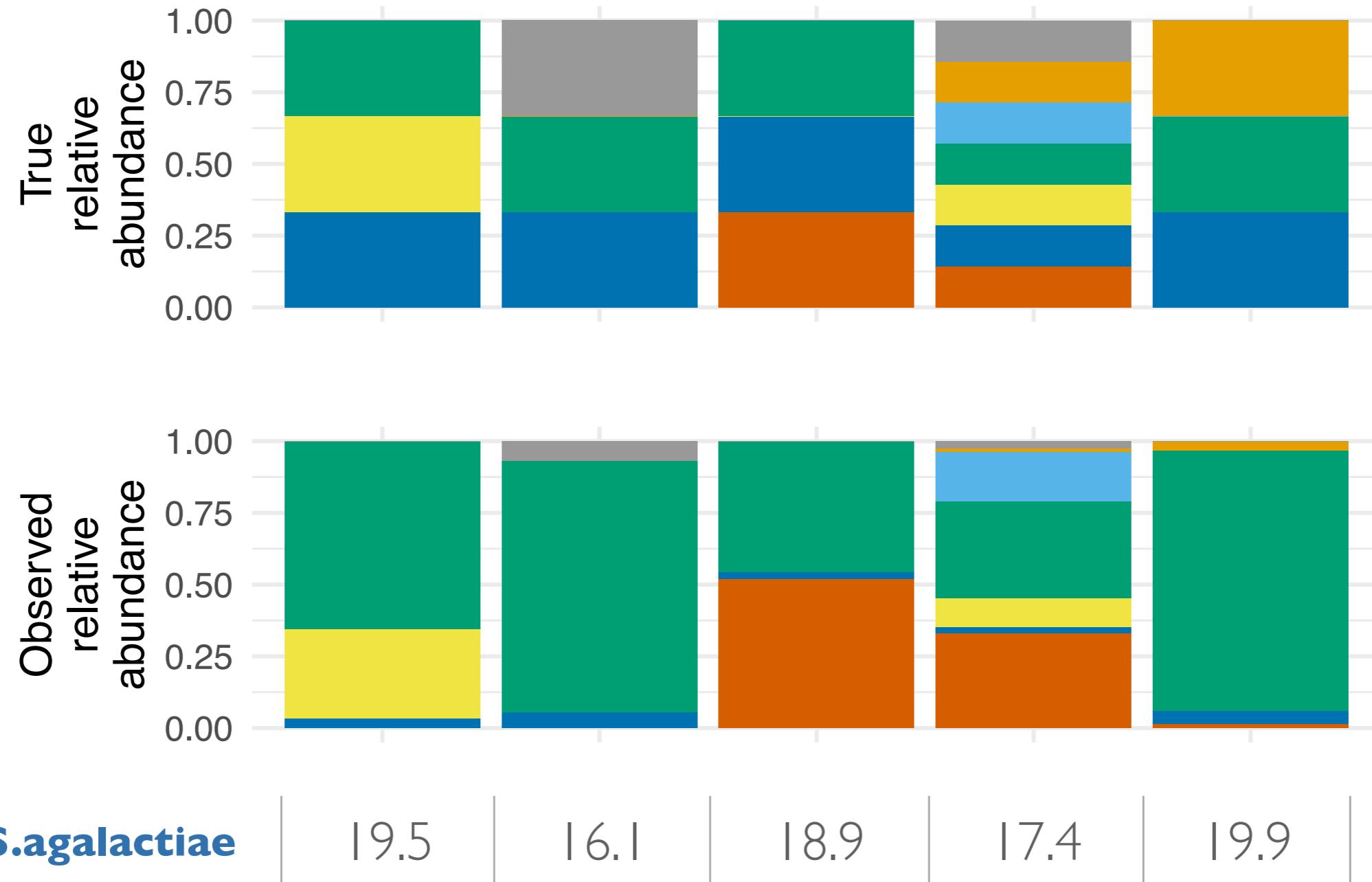
- Despite the common assumption that

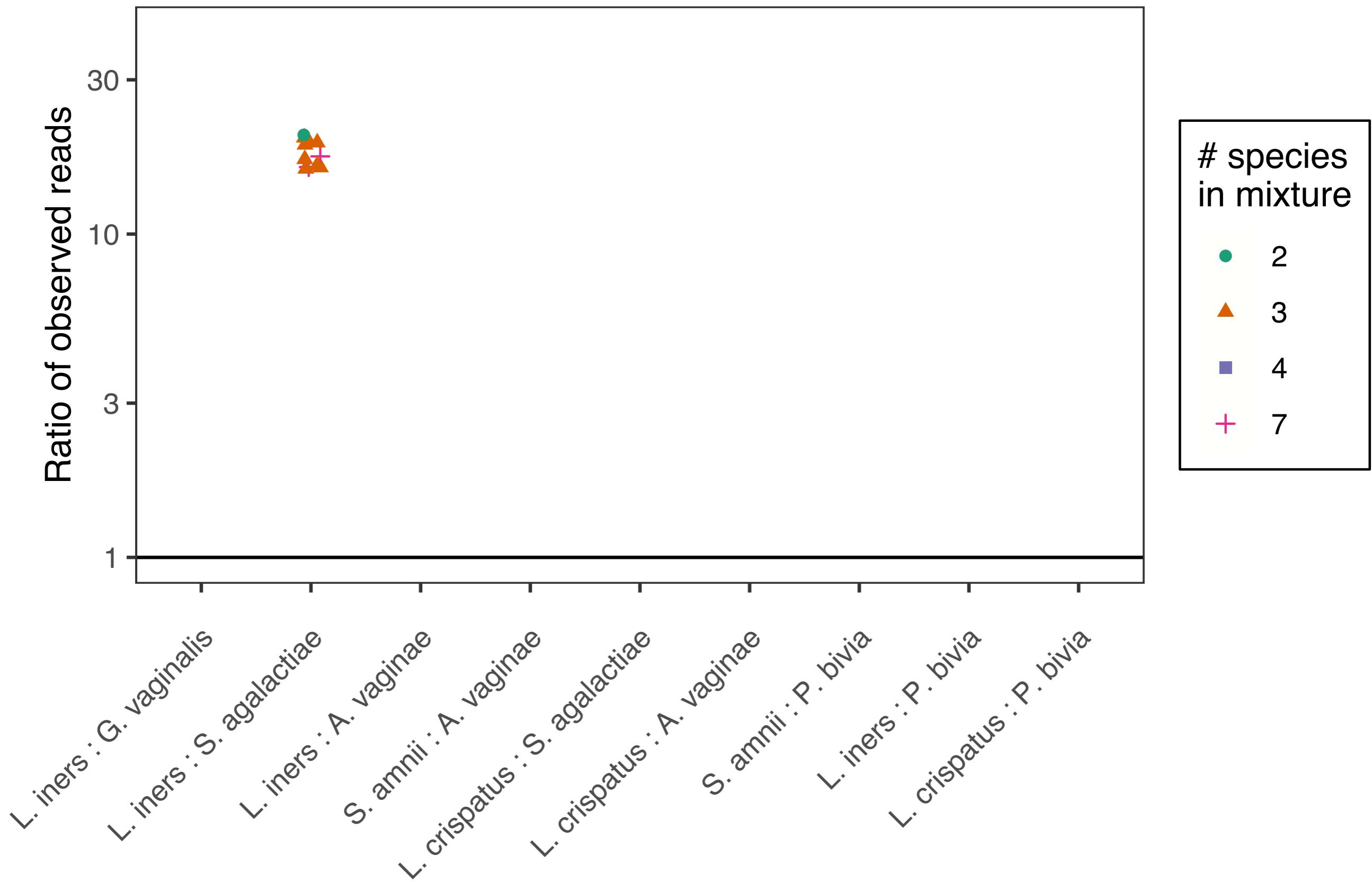
$$\mathbb{E}[W_{ij}] = c_i Y_{ij}$$

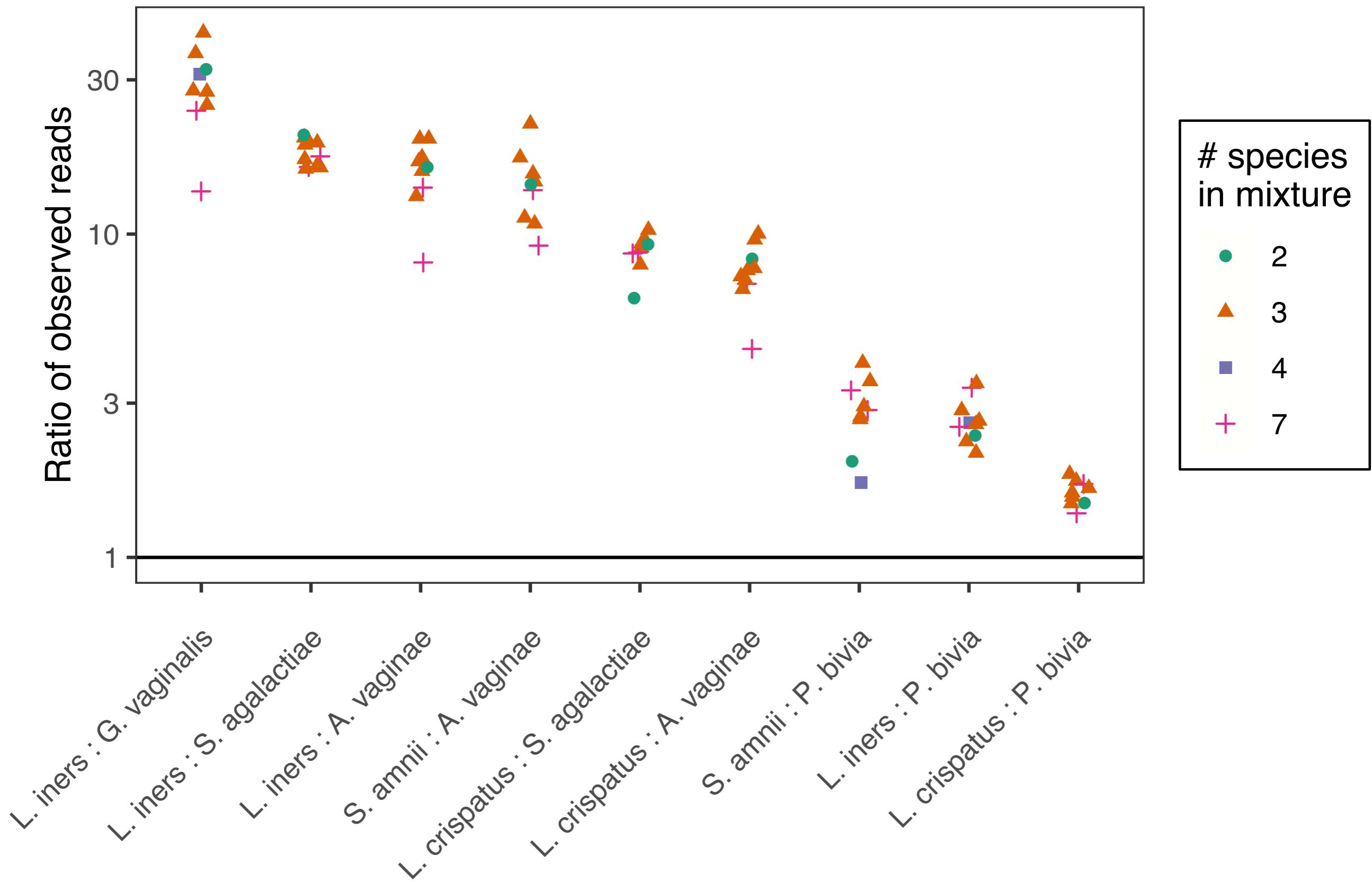
some taxa are *over-observed* for equal  $c_i$  and  $Y_{ij}$

- What model better explains this observation?









# Connecting data to reality

- Evidence *against*

$$\mathbb{E}[W_{ij}] = c_i \times Y_{ij}$$

- Better support for

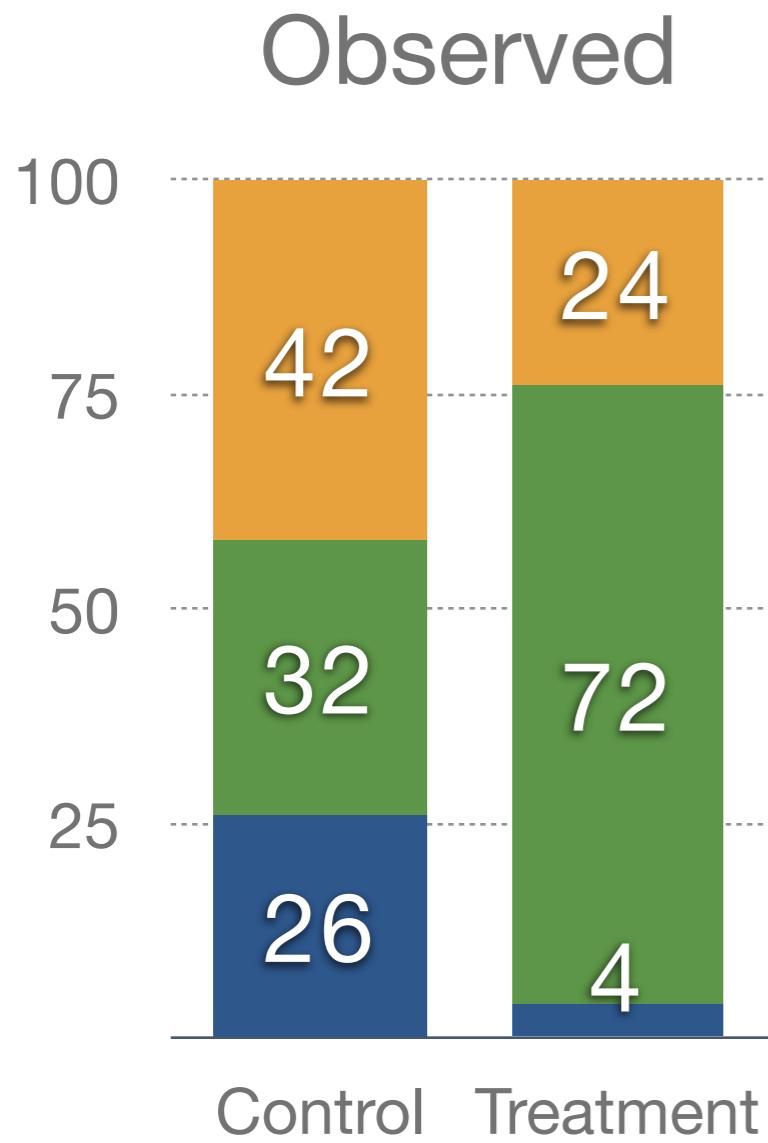
$$\mathbb{E}[W_{ij}] = c_i \times e_j \times Y_{ij}$$

Why is this so important for data analysis?

# Connecting data to reality

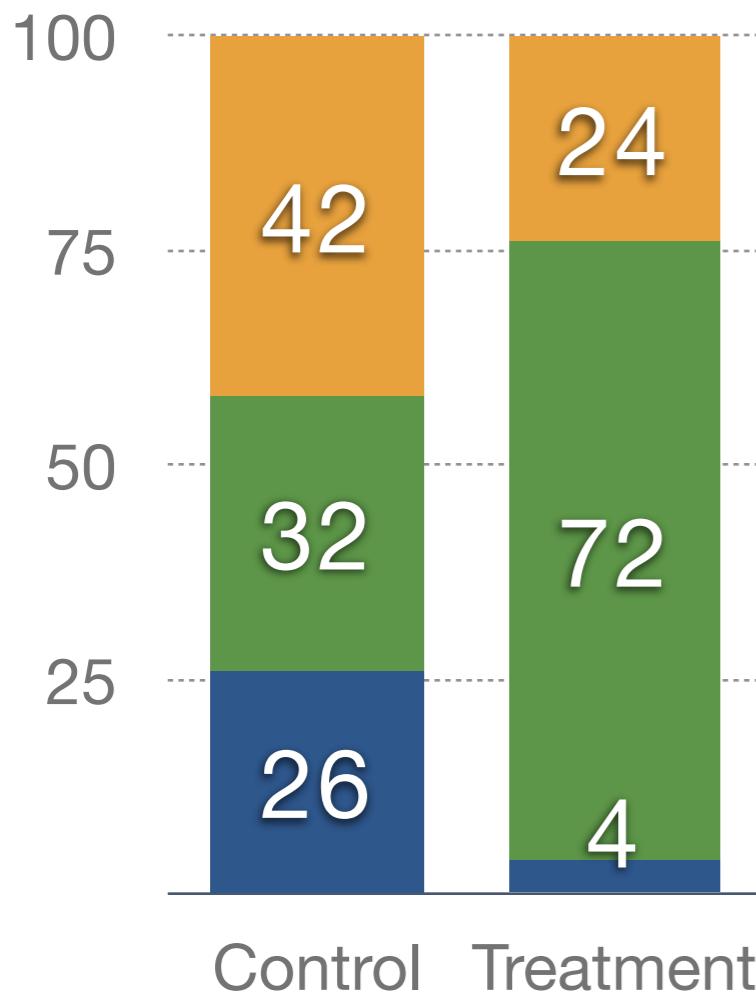
- Stated differently,

$$\text{Observed relative abundance} \propto \frac{\text{Expected value of } \frac{W_{ij}}{\sum_{j'} W_{ij'}}}{=} \frac{\text{True relative abundance} \times \text{Taxon-specific efficiencies}}{\frac{p_{ij}e_j}{\sum_{j'} p_{ij'}e_{j'}}}$$

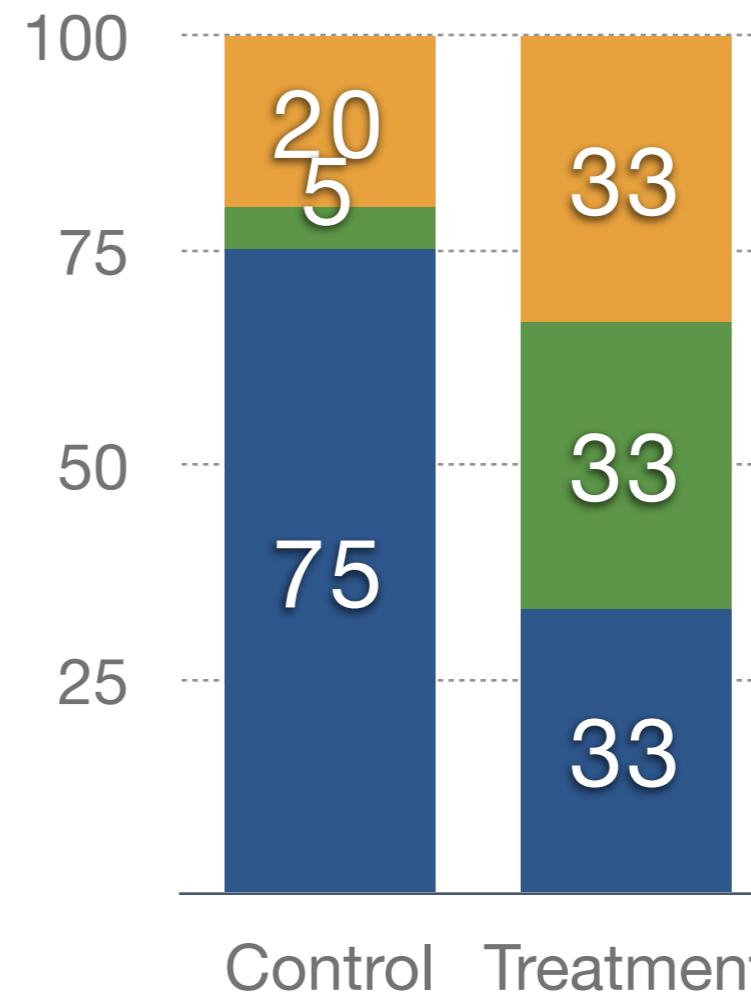



- A tempting conclusion:
  - The relative abundance of **orange** decreased in the Treatment sample (right) compared to the Control sample (left)

Observed

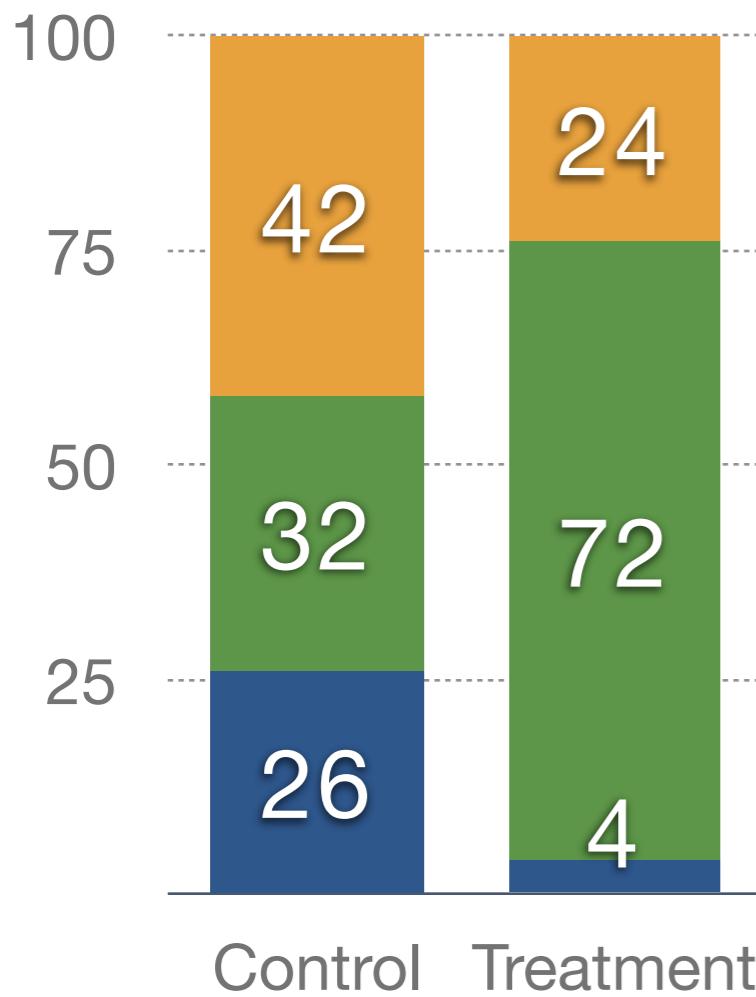


Actual

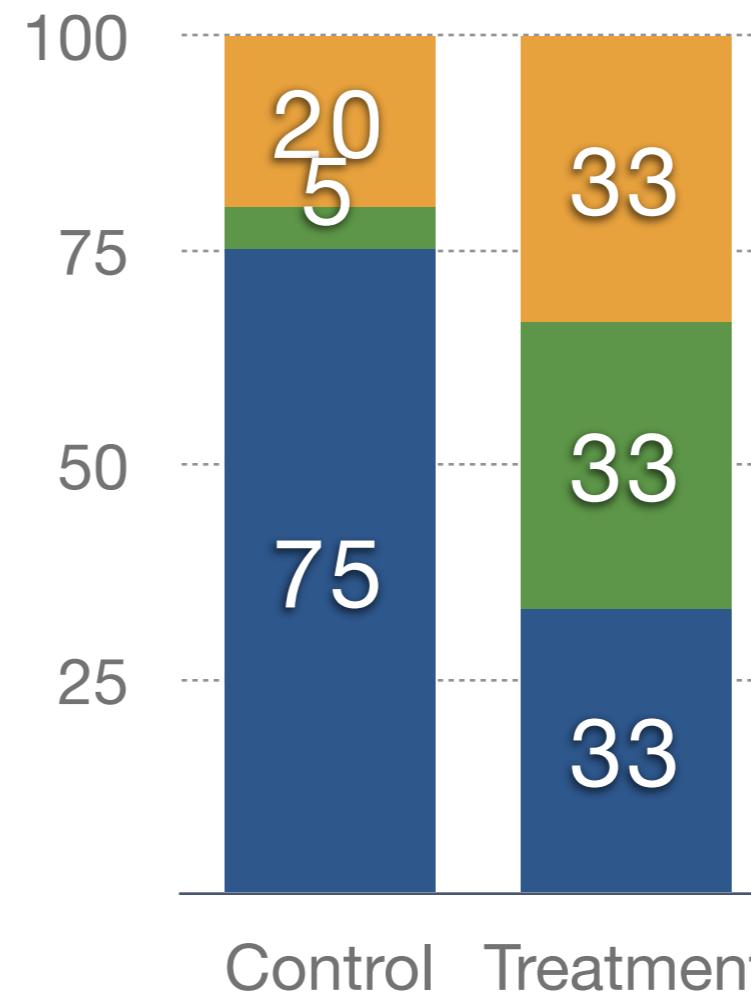


- In fact, the relative abundance of **orange increased** in the Treatment sample compared to the Control sample

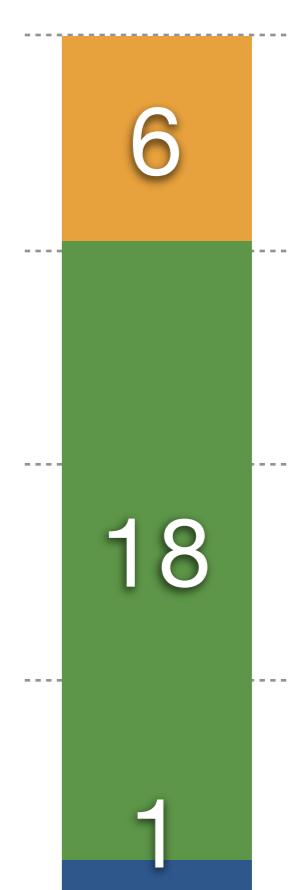
### Observed



### Actual

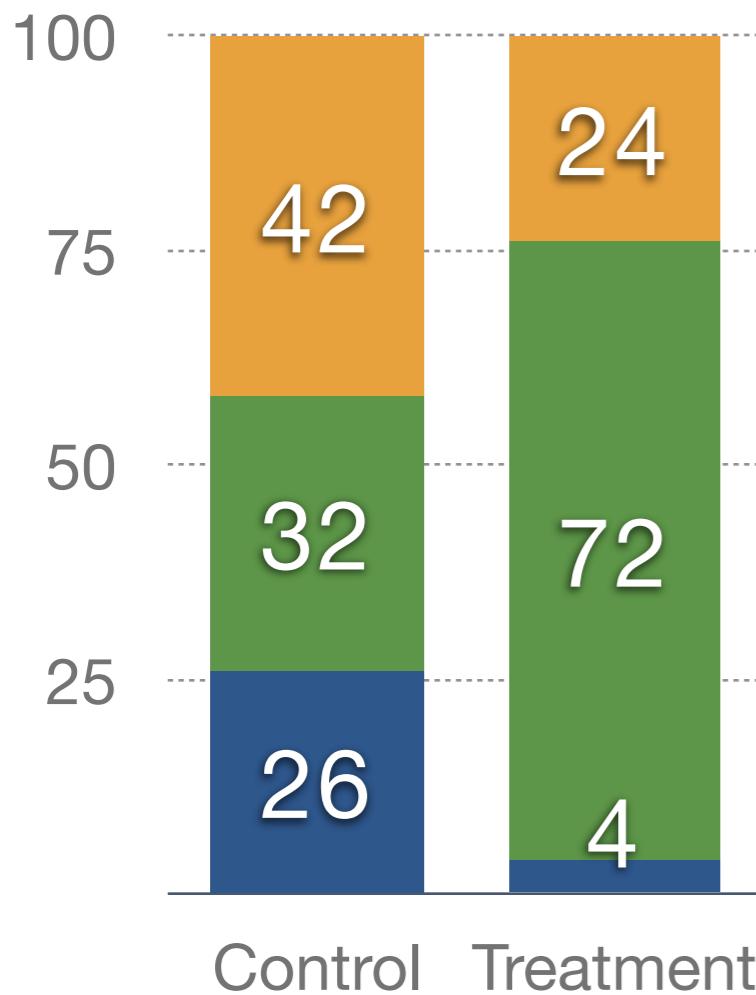


### Efficiencies

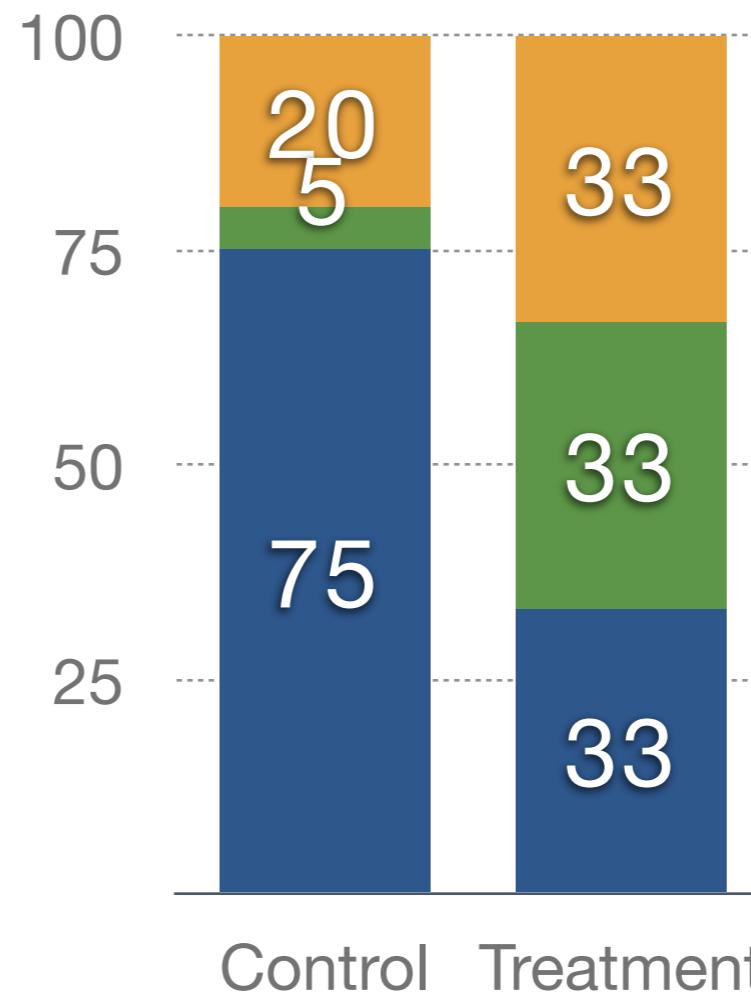


- In fact, the relative abundance of **orange increased** in the Treatment sample compared to the Control sample

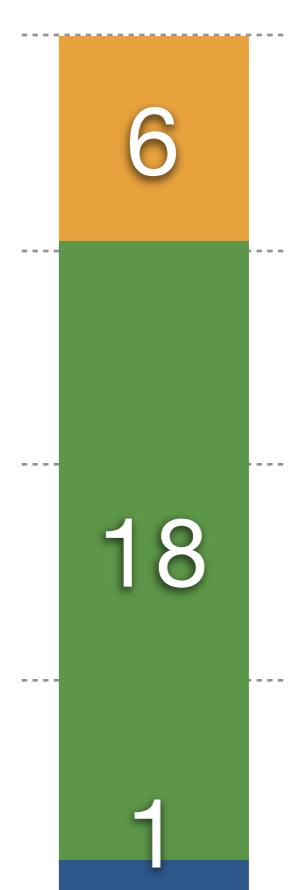
## Observed



## Actual



## Efficiencies



- **Green** is high efficiency; its abundance increased.  
**Blue** is low efficiency, and its abundance decreased.
- **Orange**'s abundance depends on the abundance of the other taxa.

# What can't we learn?

- Result: Under the model

$$\mathbb{E}W_{ij} = c_i \times e_j \times Y_{ij}$$

we cannot learn about:

- $\mathbb{E}Y_{\{X_i=1\}j} - \mathbb{E}Y_{\{X_i=0\}j}$
- $\mathbb{E}\left[\frac{Y_{\{X_i=1\}j}}{\sum_{j'} Y_{\{X_i=1\}j'}}\right] - \mathbb{E}\left[\frac{Y_{\{X_i=0\}j}}{\sum_{j'} Y_{\{X_i=0\}j'}}\right]$
- $\frac{\mathbb{E}Y_{\{X_i=1\}j}}{\mathbb{E}Y_{\{X_i=0\}j}}$  and  $\log\left(\frac{\mathbb{E}Y_{\{X_i=1\}j}}{\mathbb{E}Y_{\{X_i=0\}j}}\right)$

# What can we learn?

- Result: Under the model

$$\mathbb{E}W_{ij} = c_i \times e_j \times Y_{ij}$$

we can learn about:

- $\log\left(\frac{\mathbb{E}Y_{\{X_i=1\}j}}{\mathbb{E}Y_{\{X_i=0\}j}}\right) - \log\left(\frac{\mathbb{E}Y_{\{X_i=1\}j'}}{\mathbb{E}Y_{\{X_i=0\}j'}}\right)$
- $\log\left(\frac{\mathbb{E}Y_{\{X_i=1\}j}}{\mathbb{E}Y_{\{X_i=0\}j}}\right) - \text{average}_{j'} \log\left(\frac{\mathbb{E}Y_{\{X_i=1\}j'}}{\mathbb{E}Y_{\{X_i=0\}j'}}\right)$

# What can we learn?

- Result: Under the model

$$\mathbb{E}W_{ij} = c_i \times e_j \times Y_{ij}$$

$$\log Y_{ij} = X_i^T \beta_j$$

we can learn about

- $\beta_{kj} - \beta_{kj'}$  log ratios-of-ratios
- $\beta_{kj} - \text{average}(\beta_{k\cdot})$  log ratios *relative* to average log ratios

# What can we learn?

- Res

We can identify  
groups (taxa, genes, etc.)  
that are  
changing the **most**  
in abundance  
from HTS

we

•

•

$P_{kj}$  average ( $P_{kj}$ ) log ratios relative to average log ratios



# radEmu



- We propose an estimator of  $\beta_{kj} - \text{average}(\beta_{k\cdot})$  under the model

$$\mathbb{E}W_{ij} = c_i \times e_j \times Y_{ij}$$

$$\log Y_{ij} = X_i^T \beta_j$$

- Estimator is *consistent* under weak conditions, *efficient* under stronger conditions



# radEmu



- ✓ Estimates a ecologically-relevant, model-agnostic, interpretable parameter

"We estimate that the average abundance of *Oliverpabstia intestinalis* in metagenomes is 11 times greater after commencing dairy work, when compared to the typical fold-differences in the average abundance of taxa across these visits."

"Under the assumption that most taxa do not differ in average abundance between visits post and prior exposure, we estimate that the abundance of *Oliverpabstia intestinalis* in metagenomes from post exposure visits is 11 times greater than prior exposure visits."

- ✓ Robust to differential detection

- ✓ Adjusts for differential sequencing depth

- ✓ ...



# radEmu

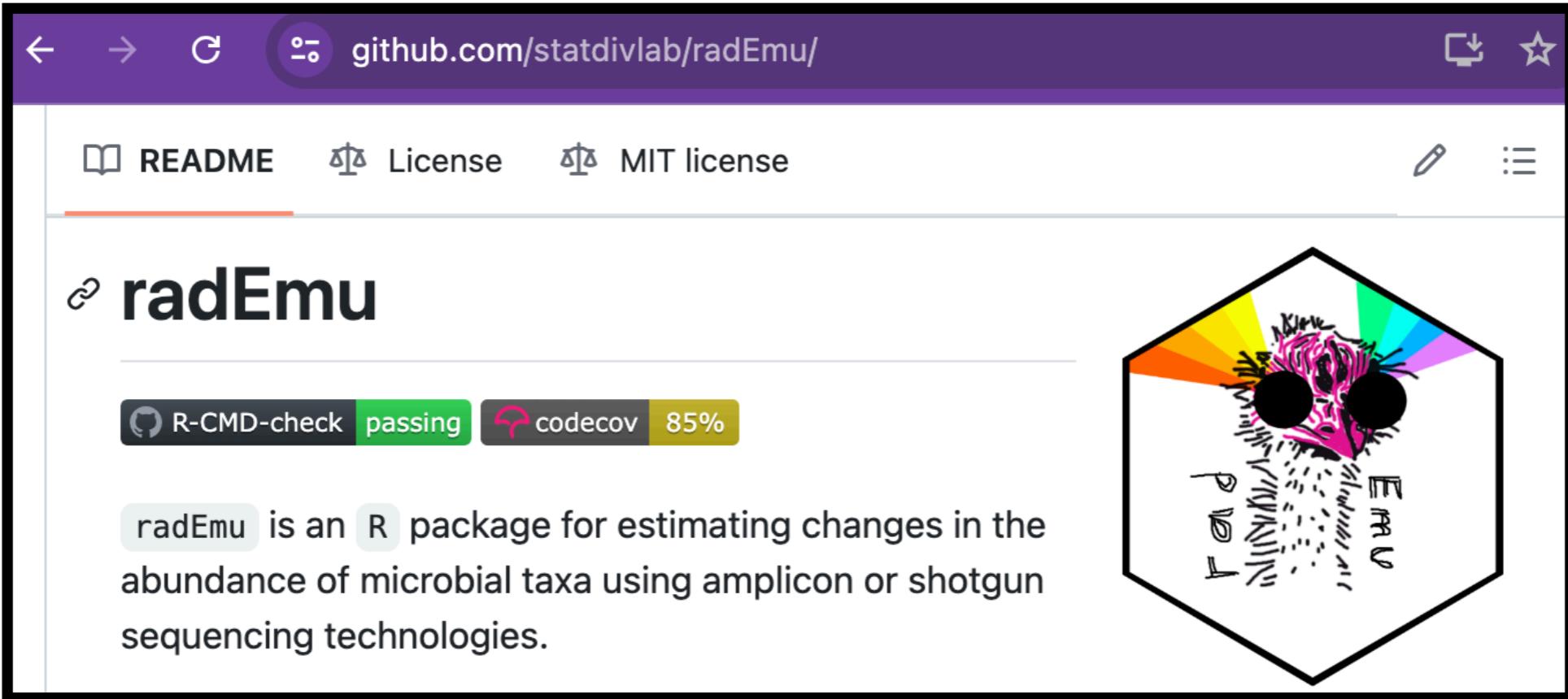


- ✓ Model-robust inference controls Type 1 error unlike DESeq2, ANCOM-BC2, t-tests...
- ✓ Handles zeroes without pseudocounts  $\mathbb{E}(\log Y_{ij})$  vs  $\log \mathbb{E}Y_{ij}$
- ✓ Simulation: Smallest estimation error out of methods that control T1E beats ALDEx2
- ✓ Covariate adjustment; inference under independence & cluster correlation;
- ✓ Spike-ins + radEmu = you can interpret fold-differences on the absolute scale
- ✗ Slower than other methods

# Summary

- “~~How do we model the data?~~”
- “How do we *learn about biology*? ”
  - Want: estimate  $\beta$  in  $\log Y_{ij} = X_i^T \beta_j$
  - Have: distorted data  $W_{ij} \approx c_i \times e_j \times Y_{ij}$
  - Result:  $\beta_{kj} - \text{average}(\beta_{k\cdot})$  is identifiable
  - Method: model-robust estimation & inference with 

# Software



The screenshot shows the GitHub repository page for `radEmu`. The URL in the address bar is `github.com/statdivlab/radEmu/`. The page features a navigation bar with links to `README`, `License`, and `MIT license`. Below this, there's a section titled `radEmu` with a brief description: "radEmu is an R package for estimating changes in the abundance of microbial taxa using amplicon or shotgun sequencing technologies." To the right of the text is a hexagonal logo depicting a stylized microorganism with internal structures and external appendages, set against a background of colored segments (yellow, green, blue, purple). At the bottom left of the page, there are status indicators for `R-CMD-check` (passing) and `codecov` (85%).

```
emuFit(formula = ~ cases + age + sex,  
        data = my_metadata,  
        Y = my_counts)
```

Statistics > Methodology

[Submitted on 7 Feb 2024 (v1), last revised 14 Mar 2025 (this version, v2)]

# Estimating Fold Changes from Partially Observed Outcomes with Applications in Microbial Metagenomics

David S Clausen, Sarah Teichman, Amy D Willis

We consider the problem of estimating fold-changes in the expected value of a multivariate outcome observed with unknown sample-specific and category-specific perturbations. This challenge arises in high-throughput sequencing studies of the abundance of microbial taxa because microbes are systematically over- and under-detected relative to their true abundances. Our model admits a partially identifiable estimand, and we establish full identifiability by imposing interpretable parameter constraints. To reduce bias and guarantee



David  
Clausen



Sarah  
Teichman



Michael  
McLaren  
(MIT,  
SecureBio)



RESEARCH ARTICLE



## Consistent and correctable bias in metagenomic sequencing experiments

Michael R McLaren<sup>1</sup>, Amy D Willis<sup>2</sup>, Benjamin J Callahan<sup>1,3\*</sup>



## Implications of taxonomic bias for microbial differential-abundance analysis

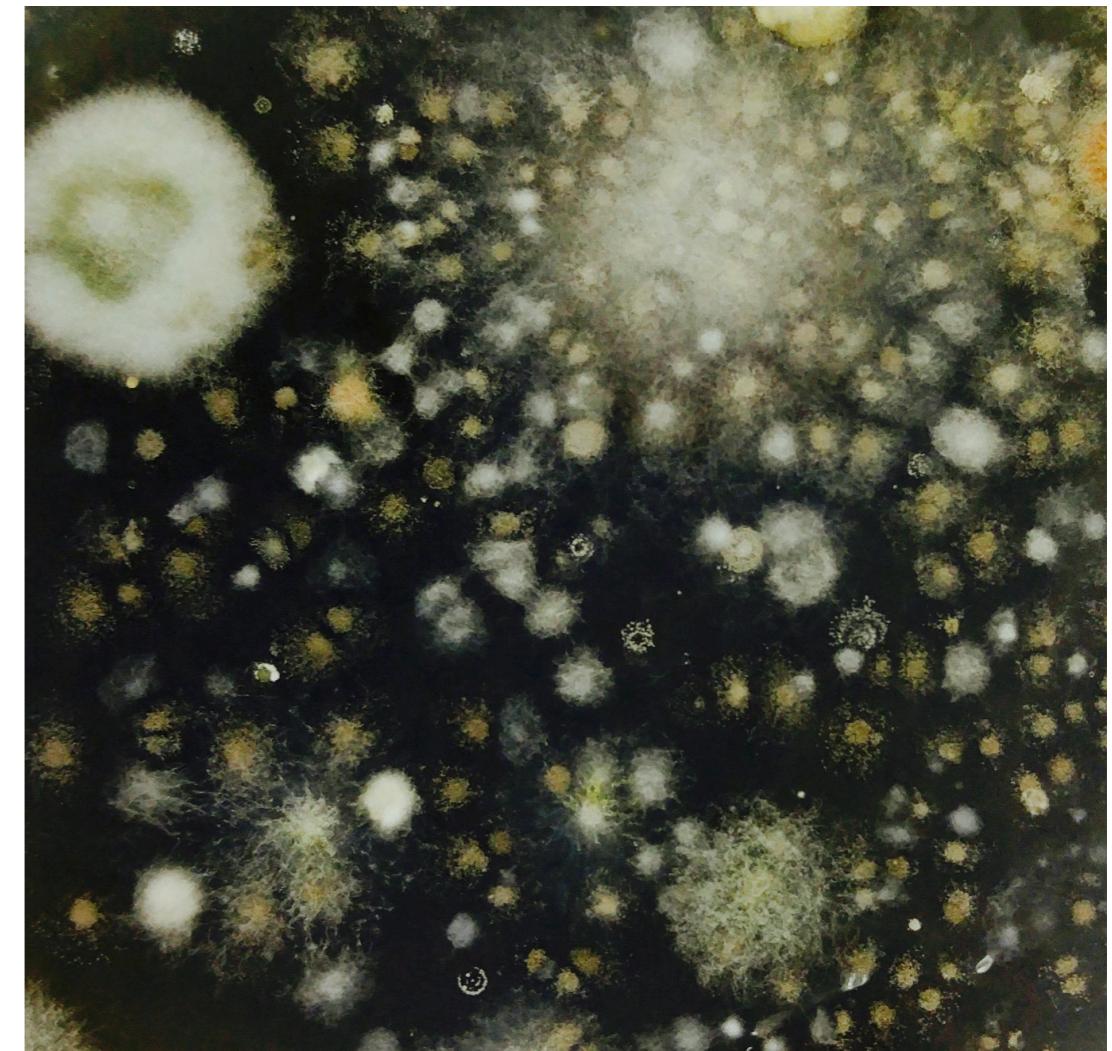
Michael R. McLaren, Jacob T. Nearing, Amy D. Willis, Karen G. Lloyd, Benjamin J. Callahan  
**doi:** <https://doi.org/10.1101/2022.08.19.504330>

Ben  
Callahan  
(NCSU)



A rigorous & rational approach to

# Microbial differential abundance



Amy D Willis PhD

Associate Professor

Department of Biostatistics

University of Washington

*Pronouns: she/her*

@AmyDWillis

[adwillis@uw.edu](mailto:adwillis@uw.edu)

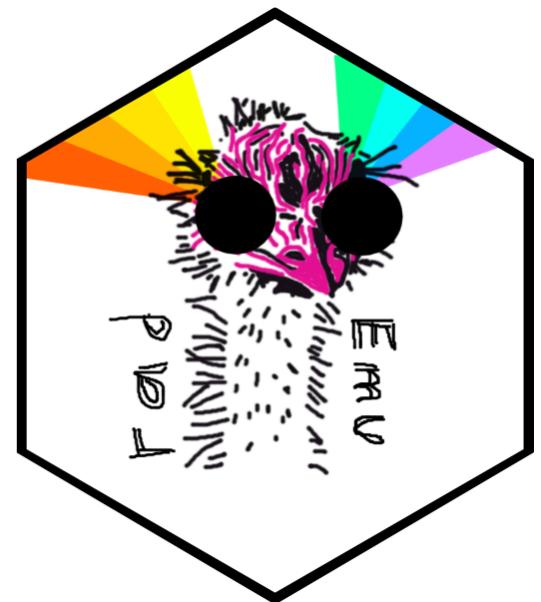


Slides: [github.com/statdivlab/presentations](https://github.com/statdivlab/presentations)

# Consistency of efficiencies

Strain	Genome size (Mbp)	Copy number	Estimated efficiency
<i>L. crispatus</i>	2.04	4	<b>2.03</b>
<i>L. iners</i>	1.30	1	<b>6.83</b>

# Amy's wish list



- You choose a meaningful parameter to estimate
- You choose a sensible way to estimate the parameter
- You choose tests that control Type 1 error

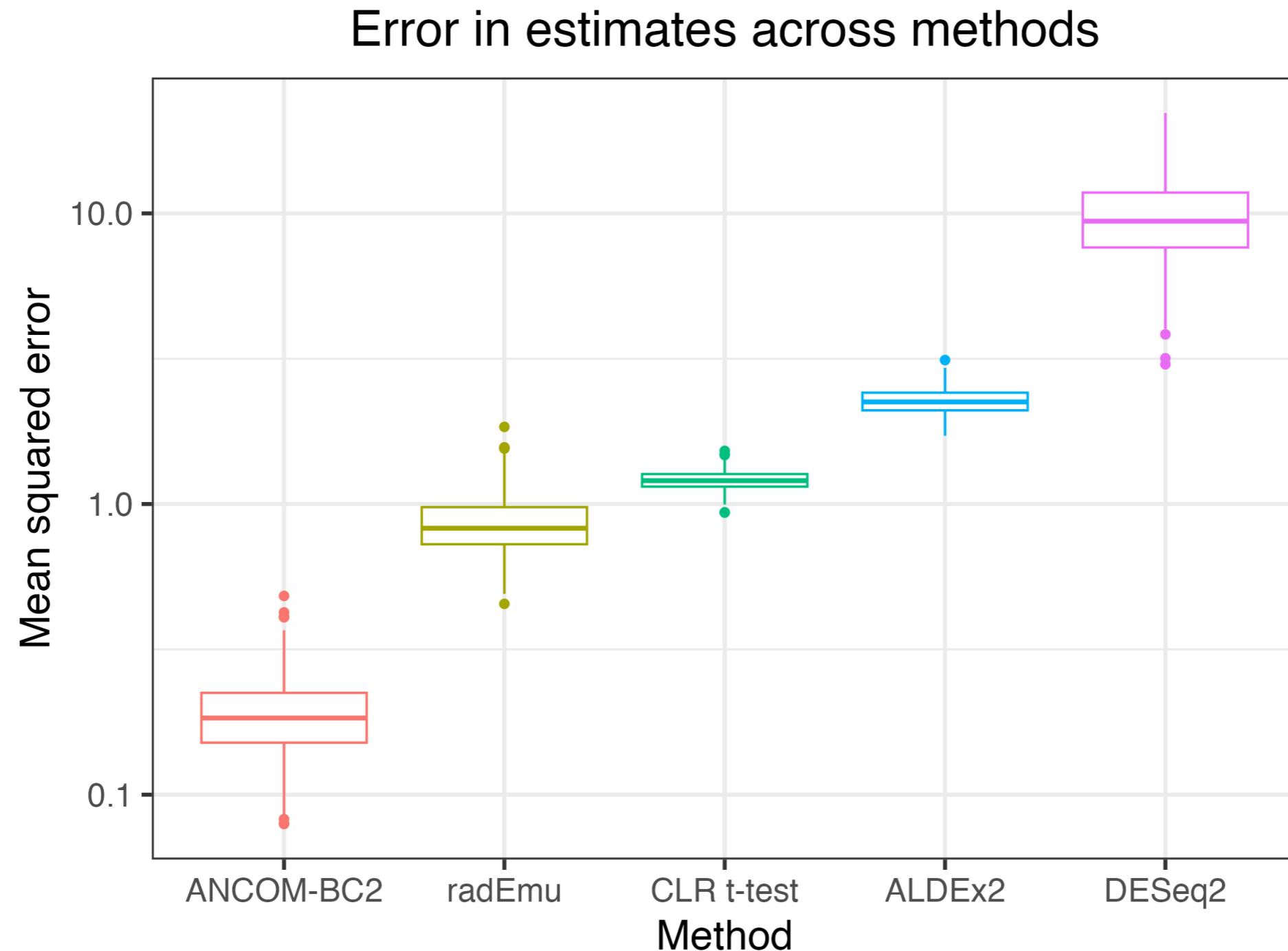
I like radEmu because it meets these criteria!

<https://github.com/statdivlab/radEmu/>

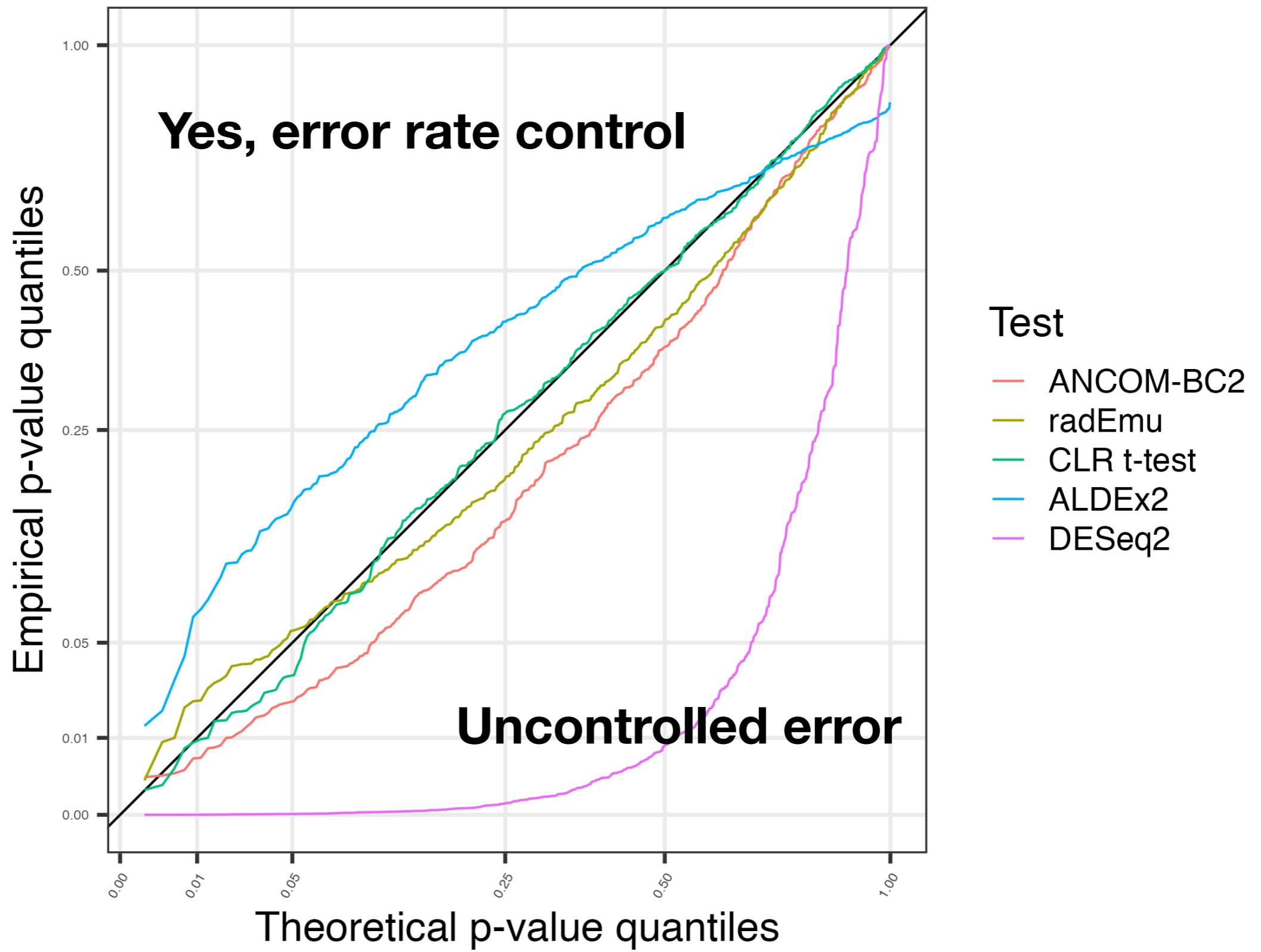
# Comparing radEmu to other methods

- Make  $W_{ij}$ 's realistic lots of zeroes, high-variance
- Ask
  1. “How good are our estimates?”
  2. “Do we have error rate control?”
    - Null hypothesis: “Fold difference (cases vs. controls) in *F. praus* is equal to typical fold difference across taxa”

# “How good are our estimates?”



# “Do we have error rate control?”



# Type I error rate control results

Method	1% Type 1 error	5% Type 1 error rate
ALDEx2	0.00	0.01
ANCOM-BC2	0.02	0.11
CLR t-test	0.01	0.06
DESeq2	0.52	0.67
radEmu	0.00	0.04

# Simulation takeaways

- TL;DR In a realistically pathological setting,
  - radEmu has the lowest error in estimation out of all methods that control error Type 1 error rate

# Identifiability: summary

- Can't estimate *absolute* log-fold differences, but can estimate differences relative to other taxa
  - $\hat{\beta}_{kj}$  is the estimated log-fold difference in the mean abundance  $\mu_{ij}$  of taxon  $j$  across  $X_k$  groups

$$\log \mathbb{E} Y_{ij}(\mathbf{x}_{x_k+1}) - \log \mathbb{E} Y_{ij}(\mathbf{x})$$

relative to typical differences across taxa

- Challenges current paradigm that we cannot make statements about *absolute abundances*  $\mu_{ij}$  from "*compositional*" data  $W_{ij}$

# Interpretation

- $\beta_{kj} - g(\beta_{k.})$  is the log-fold difference in the mean abundance of taxon  $j$  across groups that differ in  $X_{.k}$  but are similar in  $X_{\{-k\}}$ ,

$$\log \mathbb{E} Y_{ij} \left( \mathbf{x}_{x_k+1} \right) - \log \mathbb{E} Y_{ij} (\mathbf{x}),$$

relative to typical differences across taxa