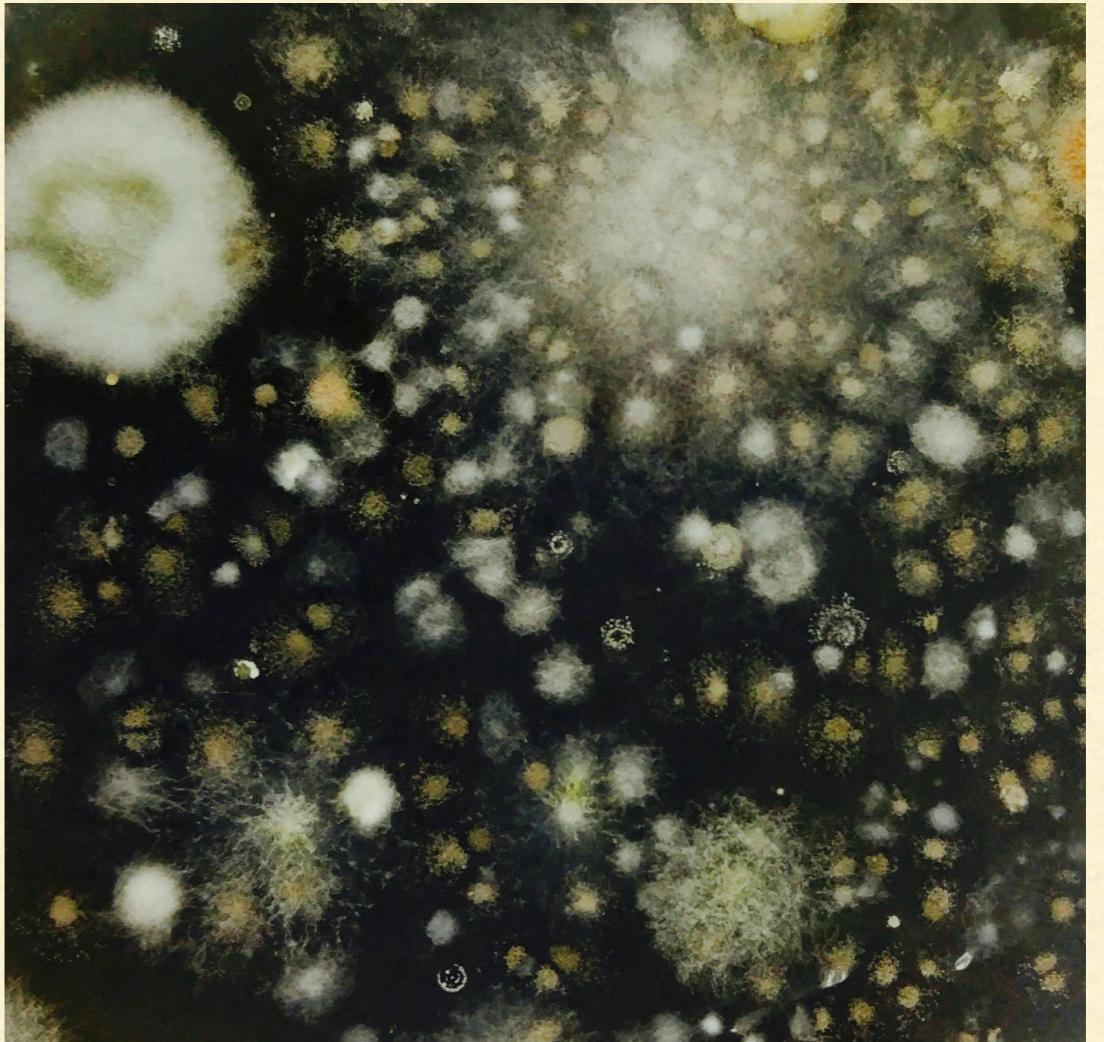


# MEASUREMENT ERROR IN MICROBIOME DATA

Amy D Willis PhD  
Assistant Professor  
Department of Biostatistics  
University of Washington, USA



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



# HIGH-THROUGHPUT SEQUENCING

---

- High-throughput sequencing (HTS) has given us high-resolution profiles of the genetic/genomic content of microbial communities
- Today's focus:

Are we *interpreting* and *analysing* HTS data correctly?

# MICROBIOME DATA FROM HTS

- At the end of the sample preparation + sequencing + bioinformatics pipeline, we often have data as counts or relative abundances

Counts

	G.vaginalis	A.vaginae	L.crispatus	L.iners
s1-1	1	1	13670	0
s1-2	1650	2930	22645	0
s1-3	1514	2073	3	1
s1-4	0	2217	16037	0
s1-5	754	0	3	0
s1-6	1	500	2	9866

Relative abundances

	G.vaginalis	A.vaginae	L.crispatus	L.iners
s1-1	0.0000	0.0000	0.6016	0.0000
s1-2	0.0606	0.1076	0.8316	0.0000
s1-3	0.0782	0.1071	0.0002	0.0001
s1-4	0.0000	0.1117	0.8081	0.0000
s1-5	0.0489	0.0000	0.0002	0.0000
s1-6	0.0001	0.0340	0.0001	0.6718

# CURRENT PARADIGM

---

- *Widely-held belief:*
  - Within a study, we can compare the relative abundance of taxa across samples
    - Directly: DeSeq2, LefSE, corncob...
    - Indirectly: Diversity (Shannon, Simpson, UniFrac), networks...



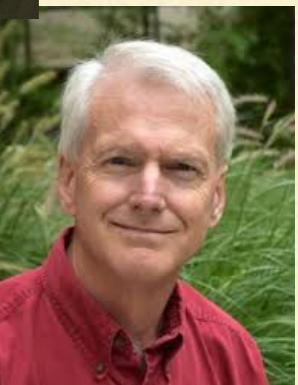
**Michael  
McLaren  
(NCSU)**



**Ben  
Callahan  
(NCSU)**



**David  
Clausen  
(UW)**



**Jim  
Hughes  
(UW)**



**Brian  
Williamson  
(UW<sup>5</sup>)**



**bioRxiv**

THE PREPRINT SERVER FOR BIOLOGY

HOME | A

Search

New Results

Comment on this paper

## A multi-view model for relative and absolute microbial abundances

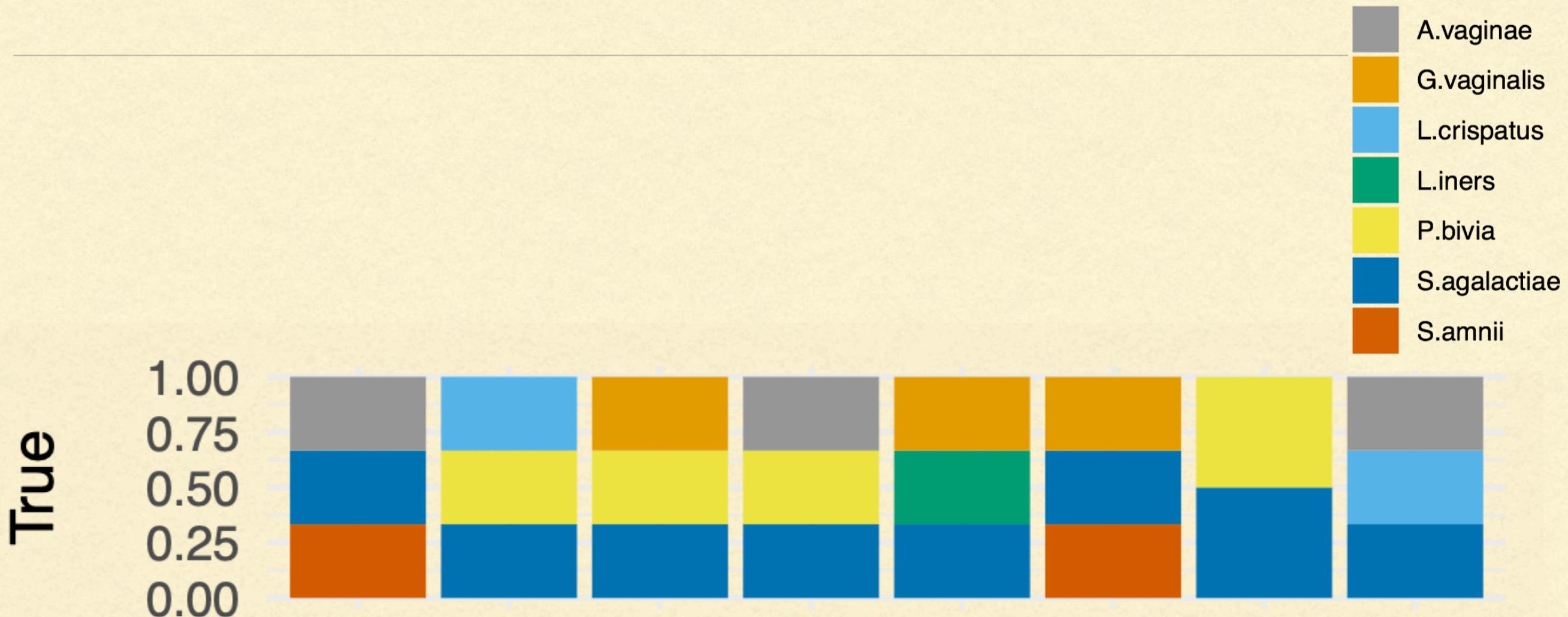
Brian D. Williamson, James P. Hughes, Amy D. Willis

doi: <https://doi.org/10.1101/761486>

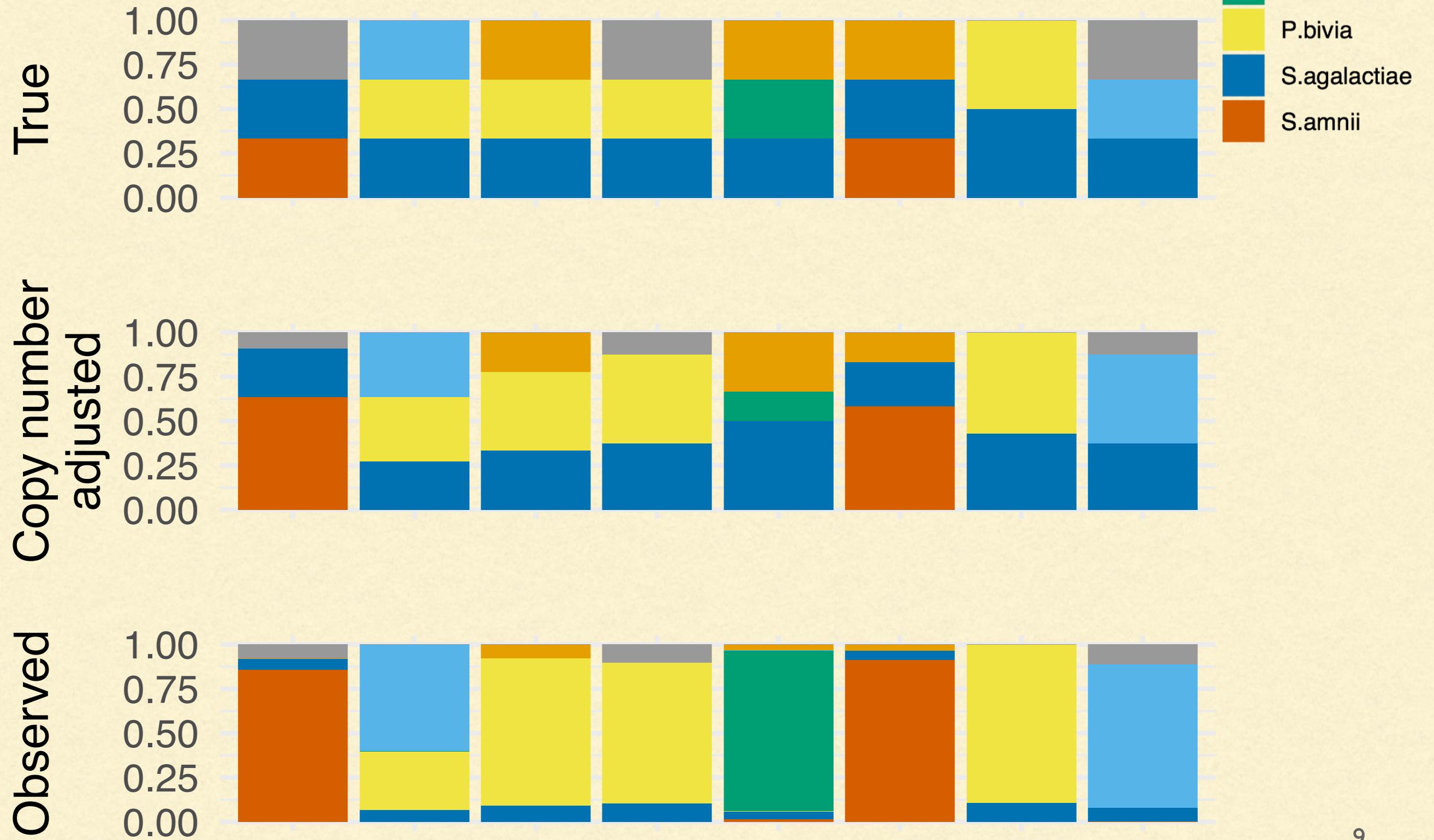
# EVALUATING RELATIVE ABUNDANCE DATA

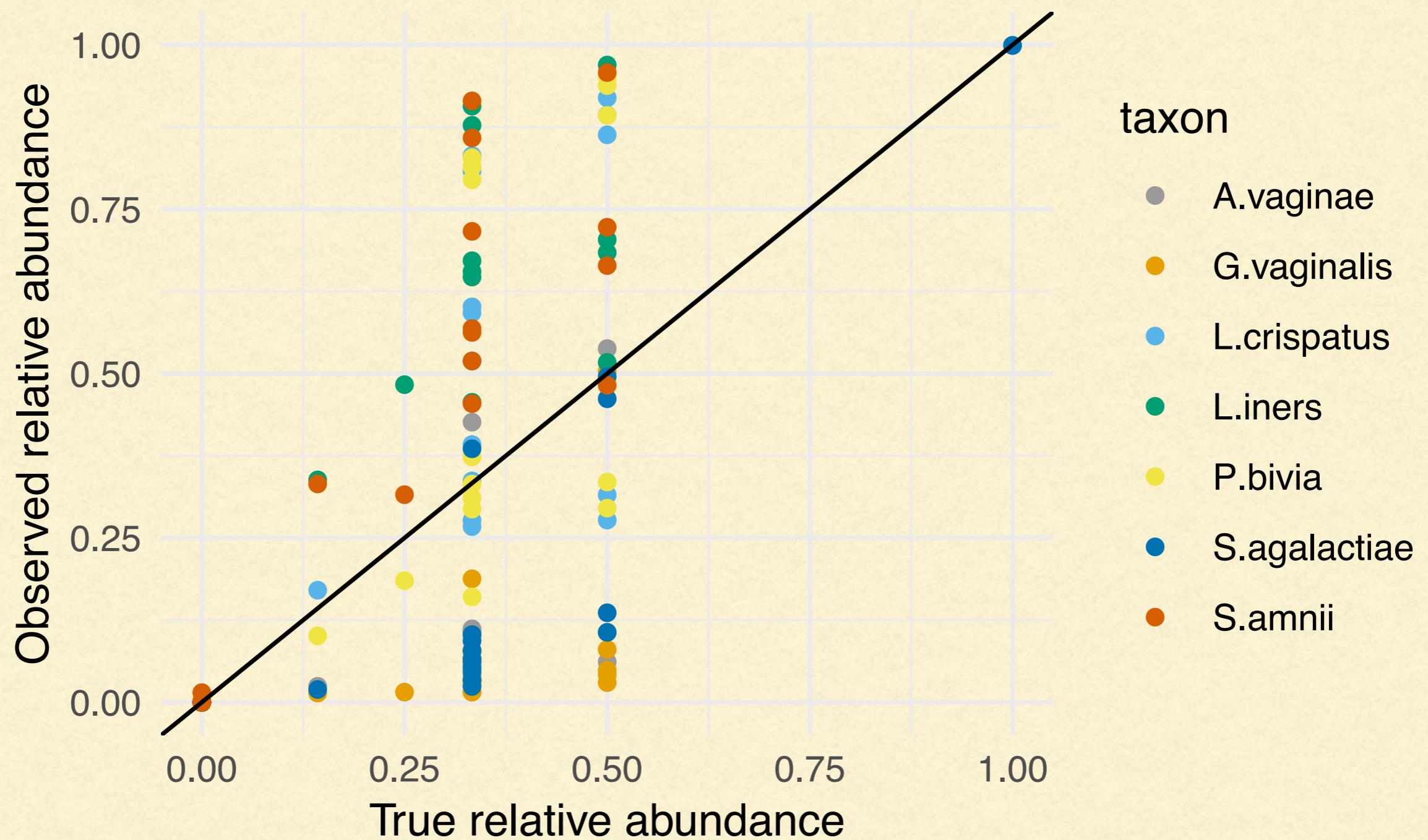
---

- Mock communities
  - Artificially constructed communities of known composition
  - Commonly used to benchmark sequencing and bioinformatics pipelines

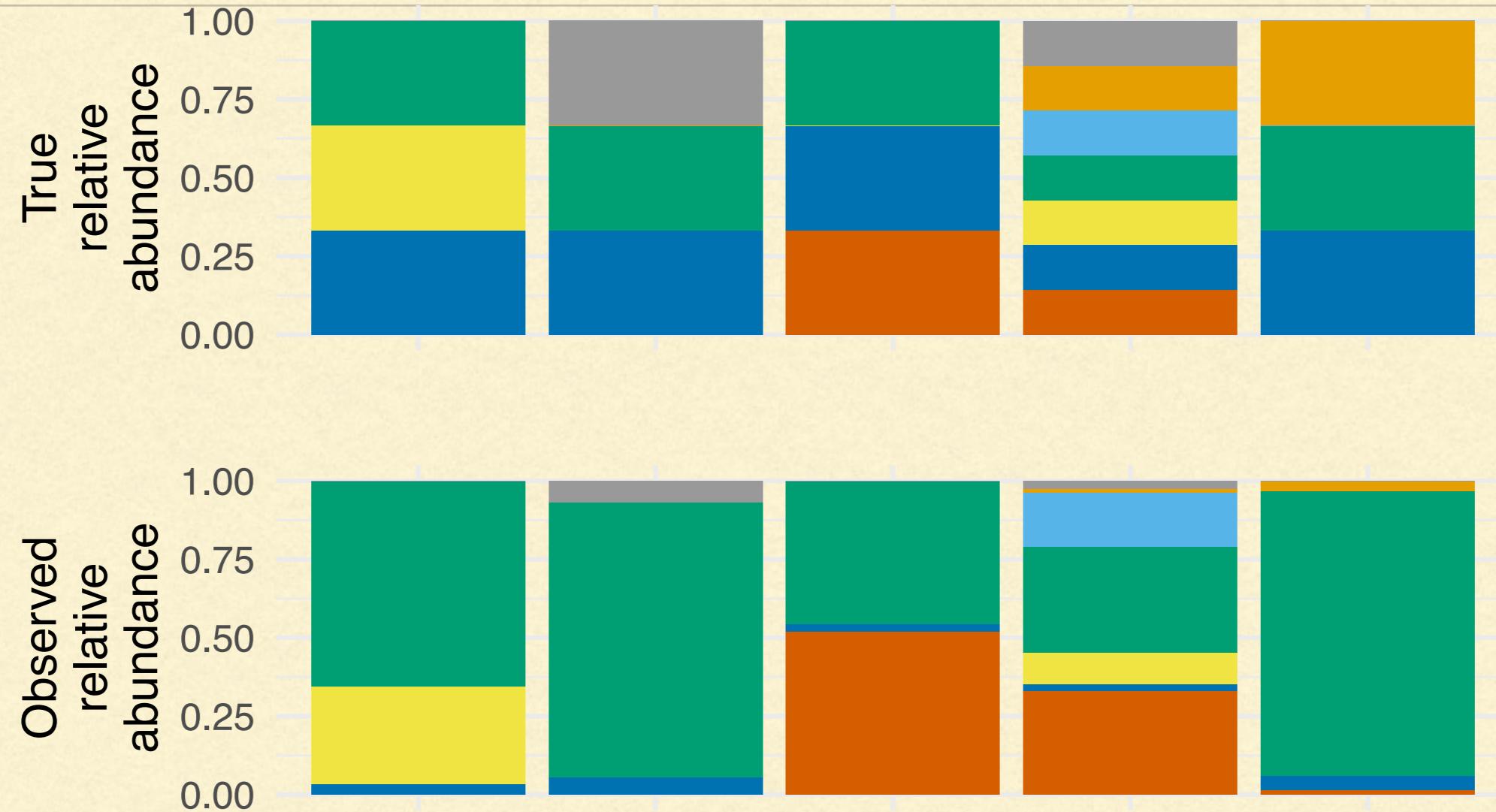




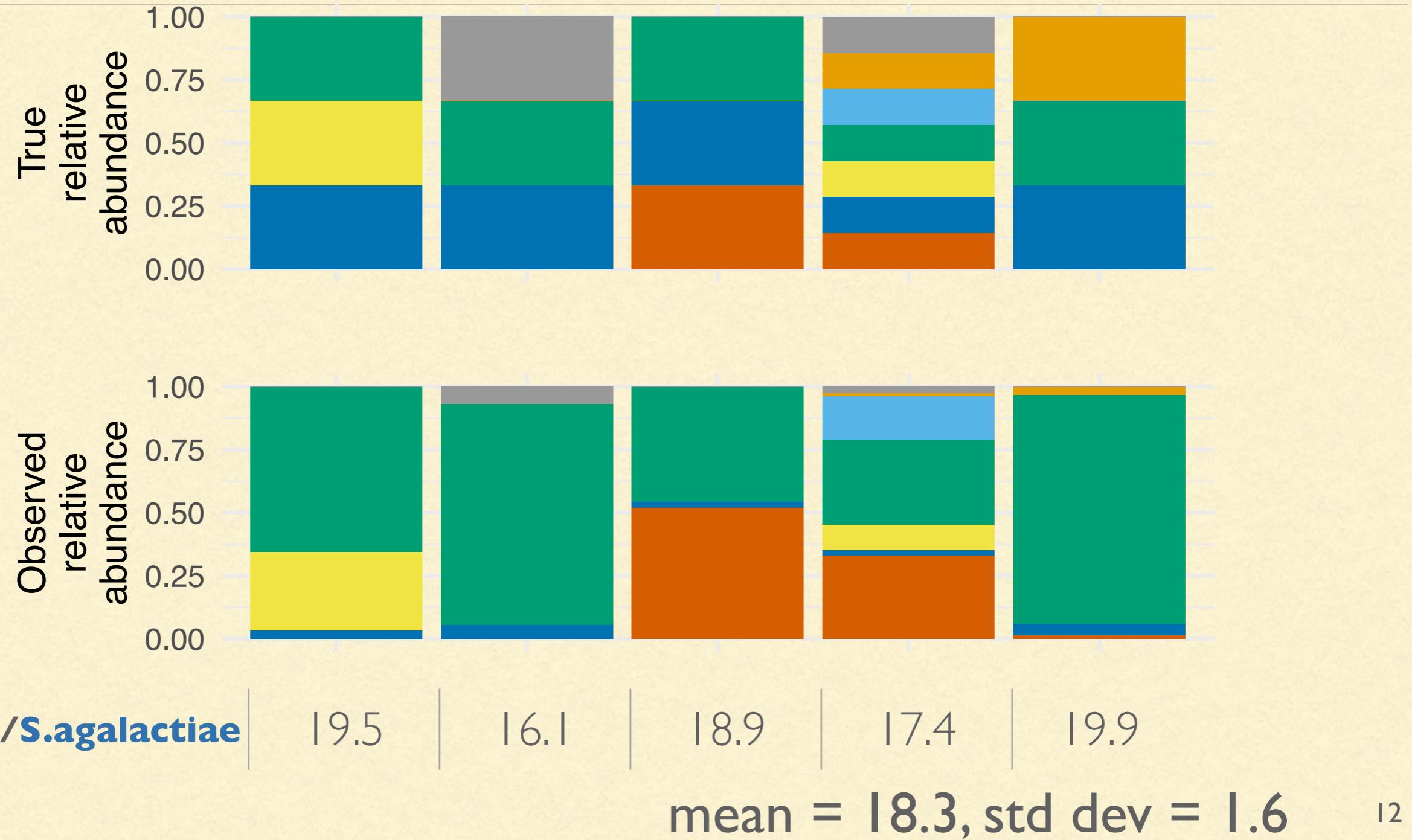


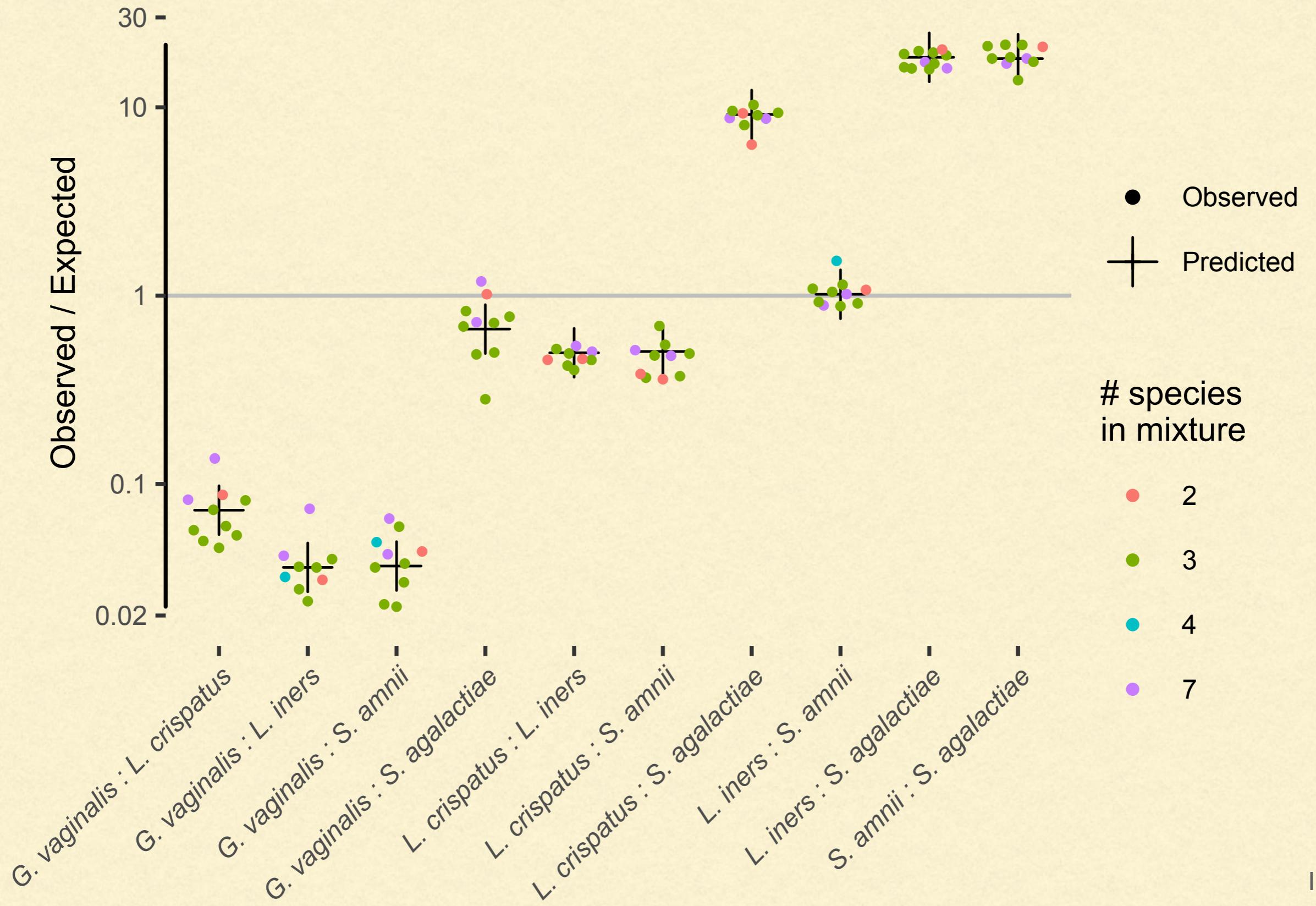


# WHAT PATTERNS CAN WE FIND?



# WHAT PATTERNS CAN WE FIND?





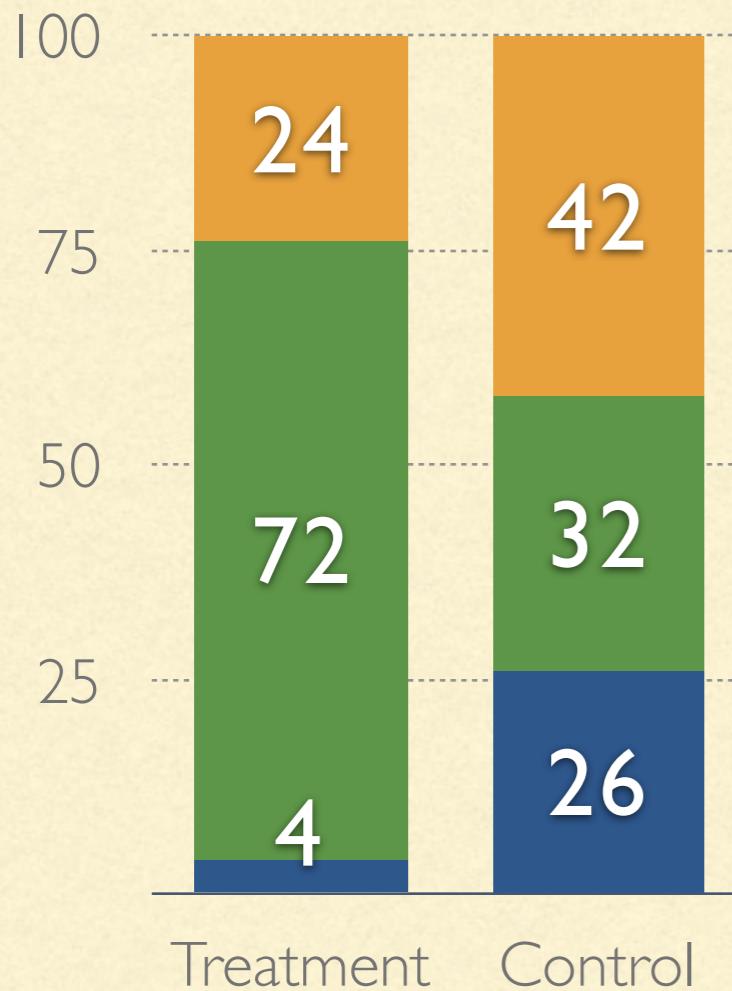
# MEASUREMENT ERROR MODEL

- This data suggests the following model:

$$\text{Observed relative abundance} \propto \text{True relative abundance} \times \text{Taxon-specific efficiencies}$$
$$\text{Expected value of } \frac{W_{ij}}{\sum_k W_{ik}} = \frac{p_{ij}e_j}{\sum_k p_{ik}e_k}$$

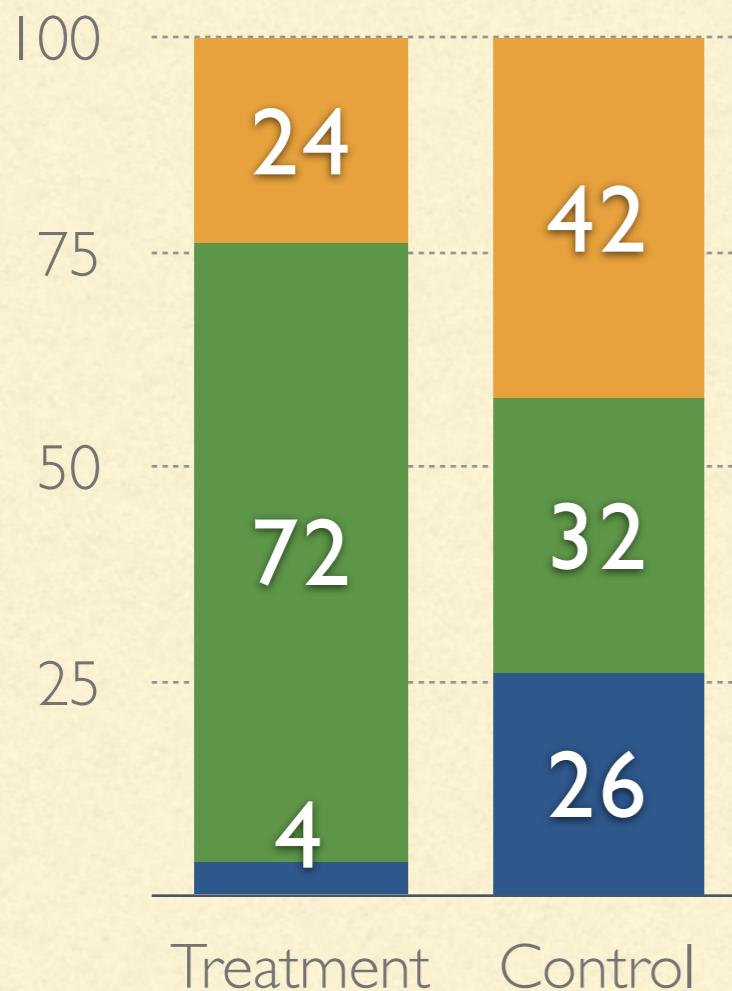
The diagram illustrates the components of the measurement error model. On the left, 'Observed relative abundance' is shown with a proportionality symbol ( $\propto$ ) above it. To its right is the mathematical expression. The expression is divided into three main parts by multiplication signs: 'True relative abundance', 'Taxon-specific efficiencies', and a fraction. Arrows point from the text labels to their corresponding parts in the expression. The 'True relative abundance' arrow points to the term  $p_{ij}e_j$ . The 'Taxon-specific efficiencies' arrow points to the term  $\sum_k p_{ik}e_k$ . The first part of the expression,  $\frac{W_{ij}}{\sum_k W_{ik}}$ , is labeled 'Expected value of' above it.

## Observed

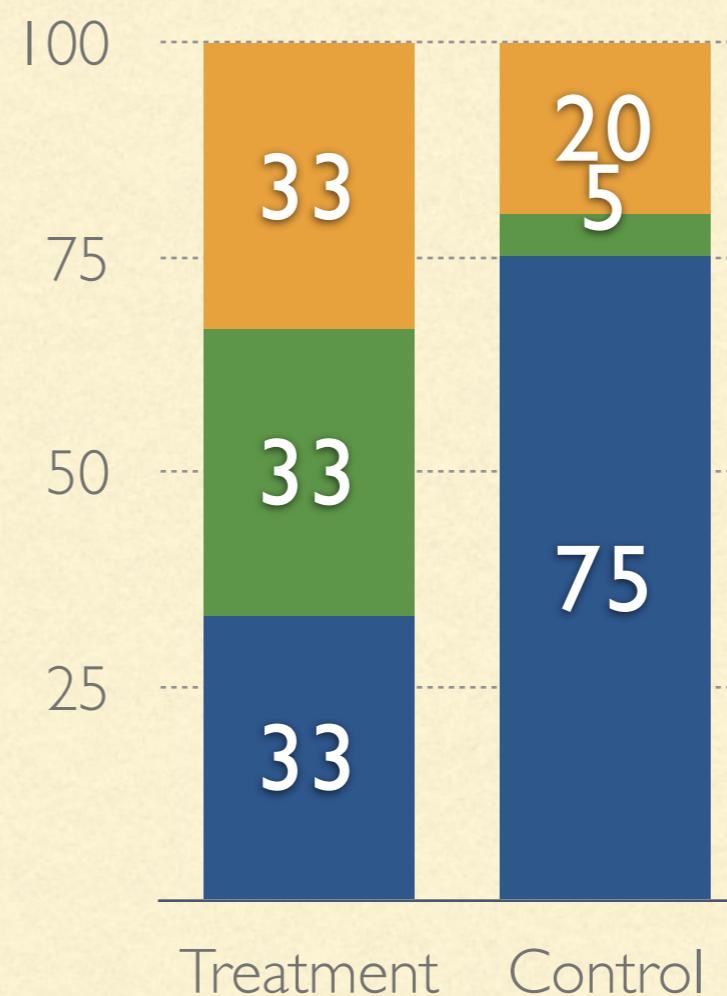


- A tempting conclusion:
  - The relative abundance of **taxon orange** decreased in the Treatment sample (left) compared to the Control sample (right)

## Observed

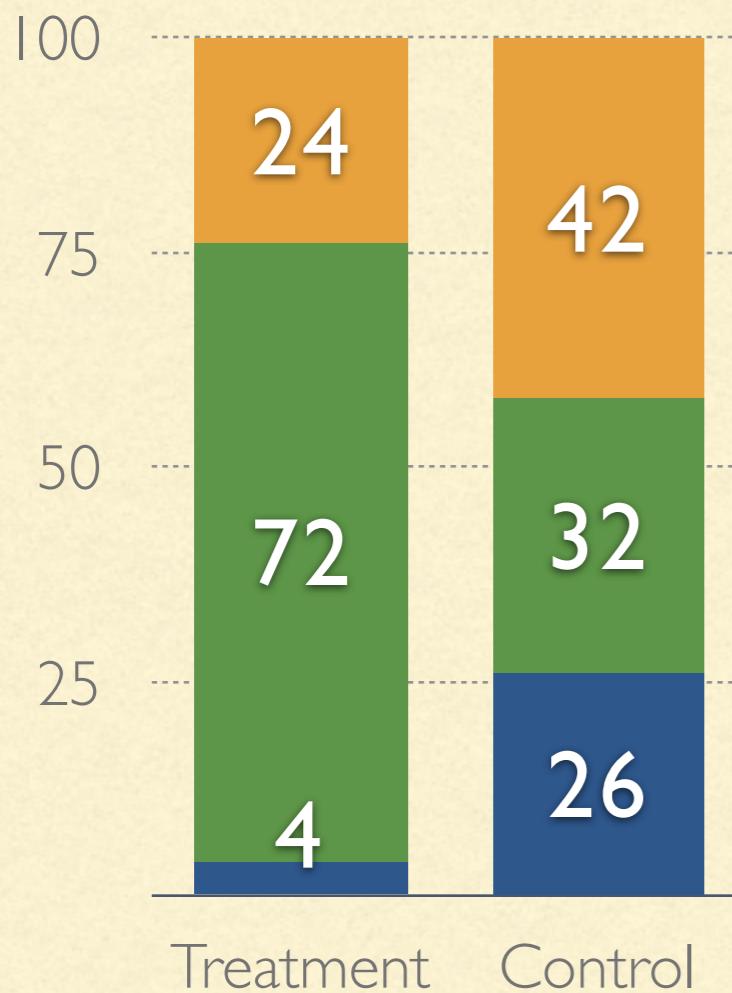


## Actual

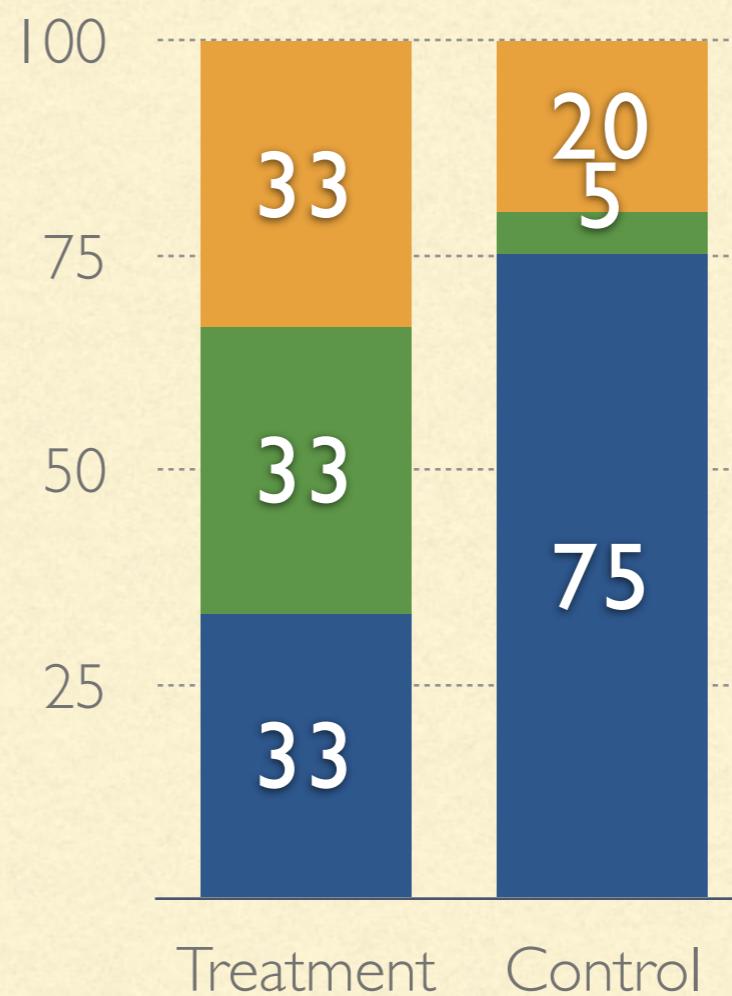


- In fact, the relative abundance of **taxon orange** increased in the Treatment sample compared to the Control sample

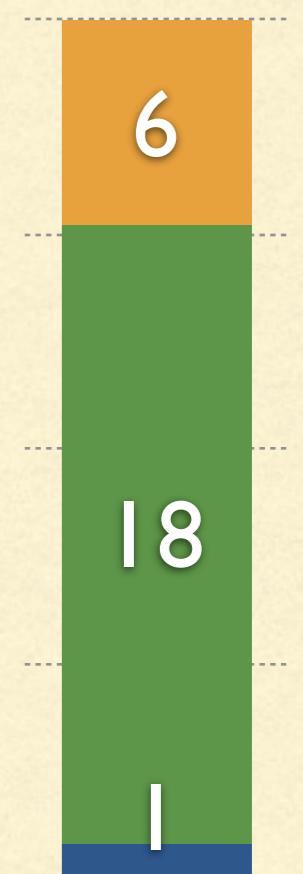
## Observed



## Actual

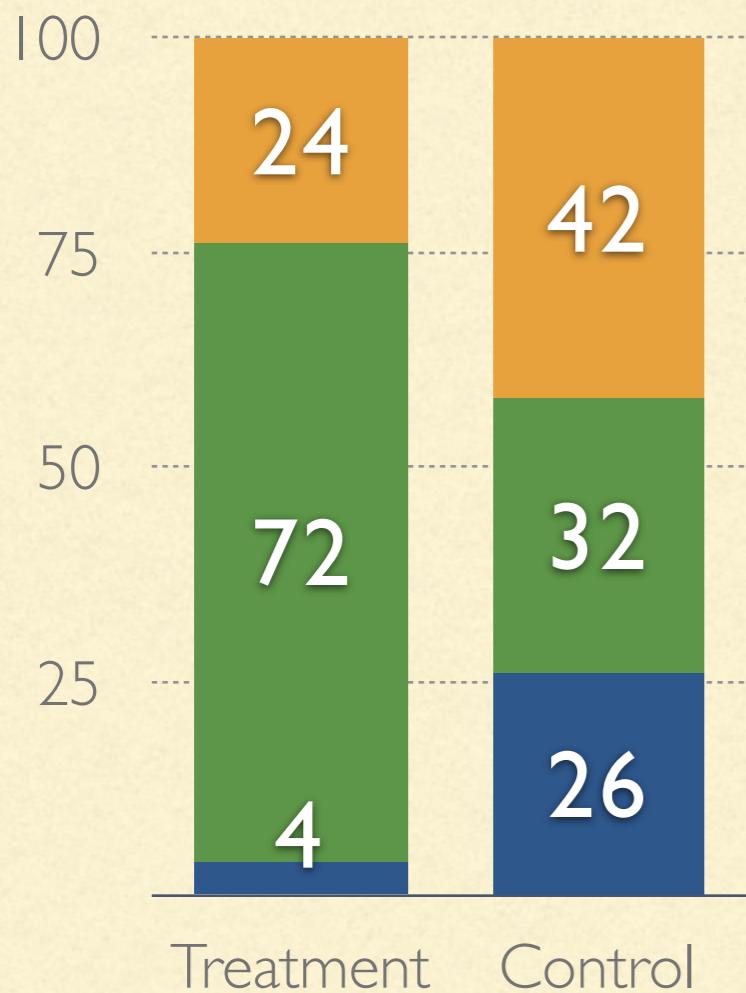


## Efficiencies

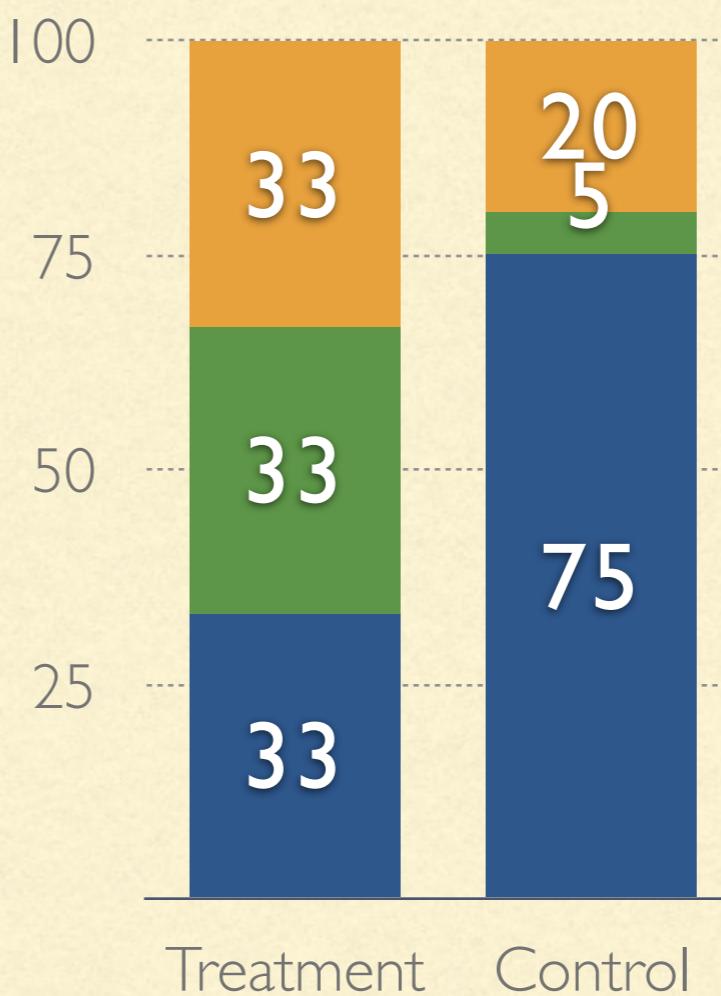


- In fact, the relative abundance of **taxon orange** increased in the Treatment sample compared to the Control sample

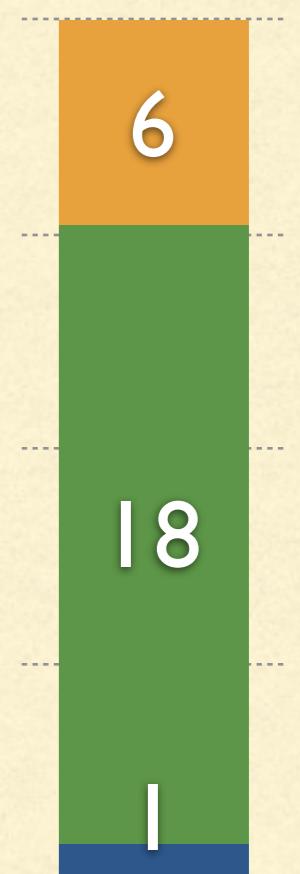
## Observed



## Actual

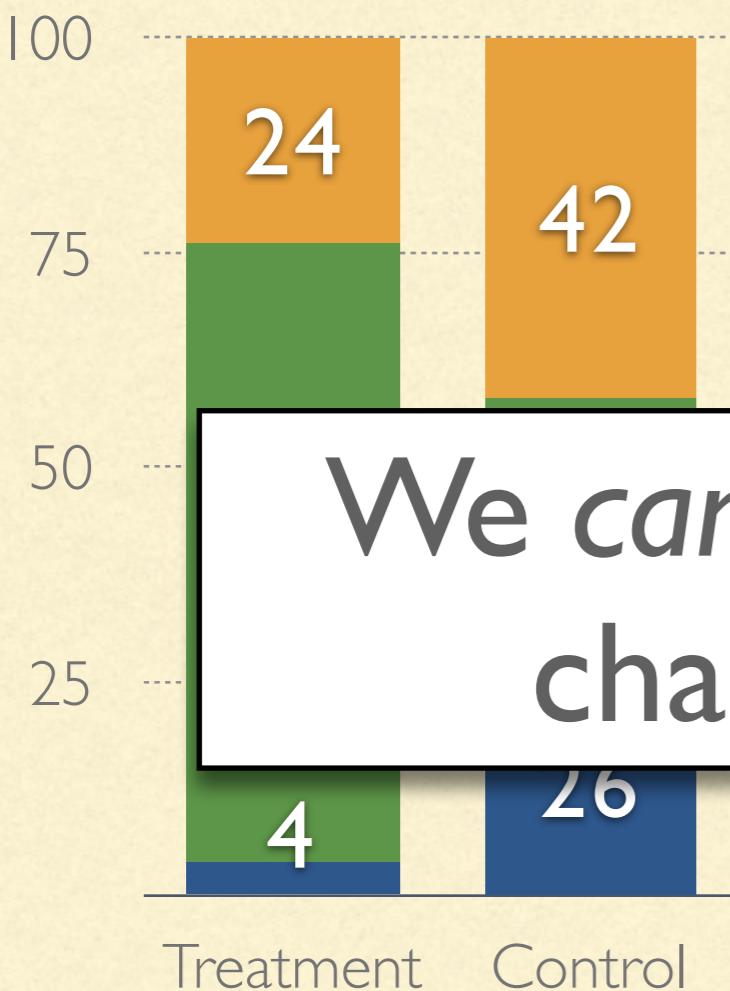


## Efficiencies

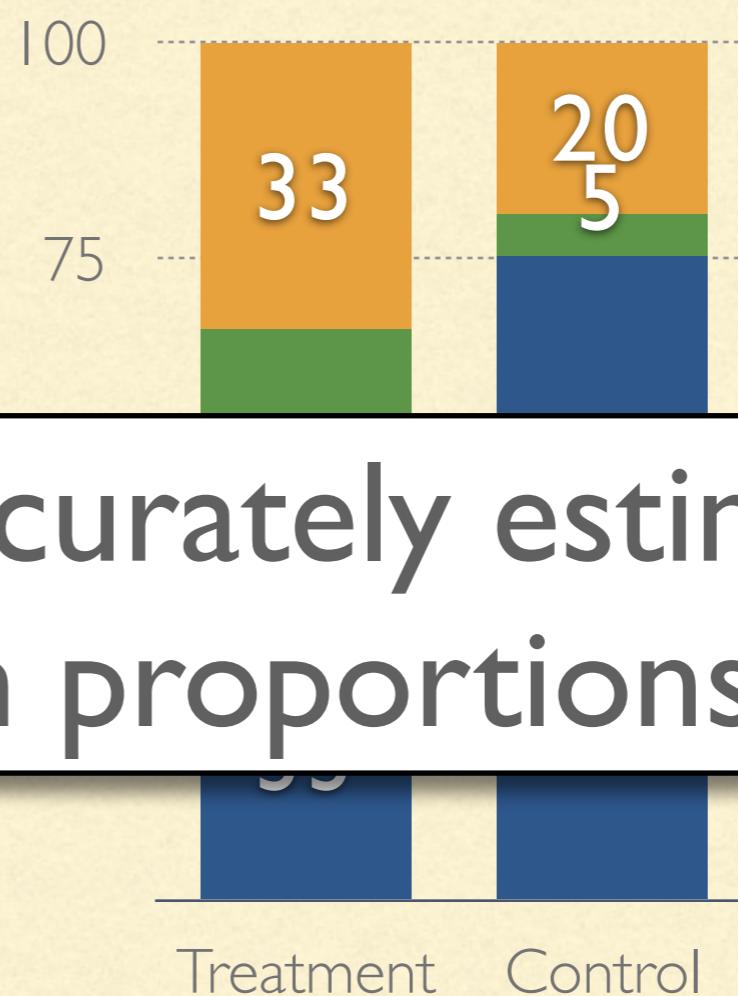


- **Taxon green** is high efficiency; its abundance increased (Ctrl vs Tmt). Additionally, **taxon blue** is low efficiency, and its abundance decreased.
- **Taxon orange**'s abundance depends on the abundance of the other taxa

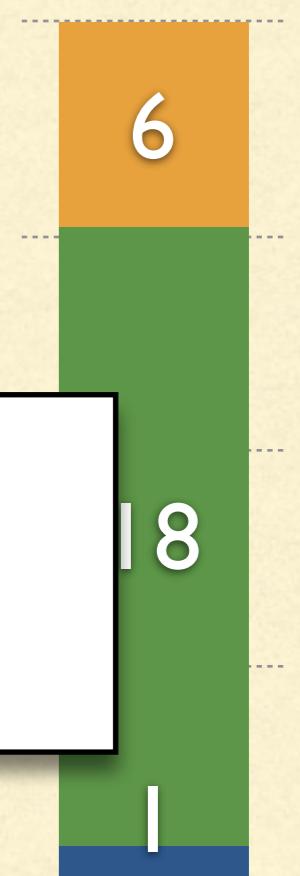
Observed



Actual



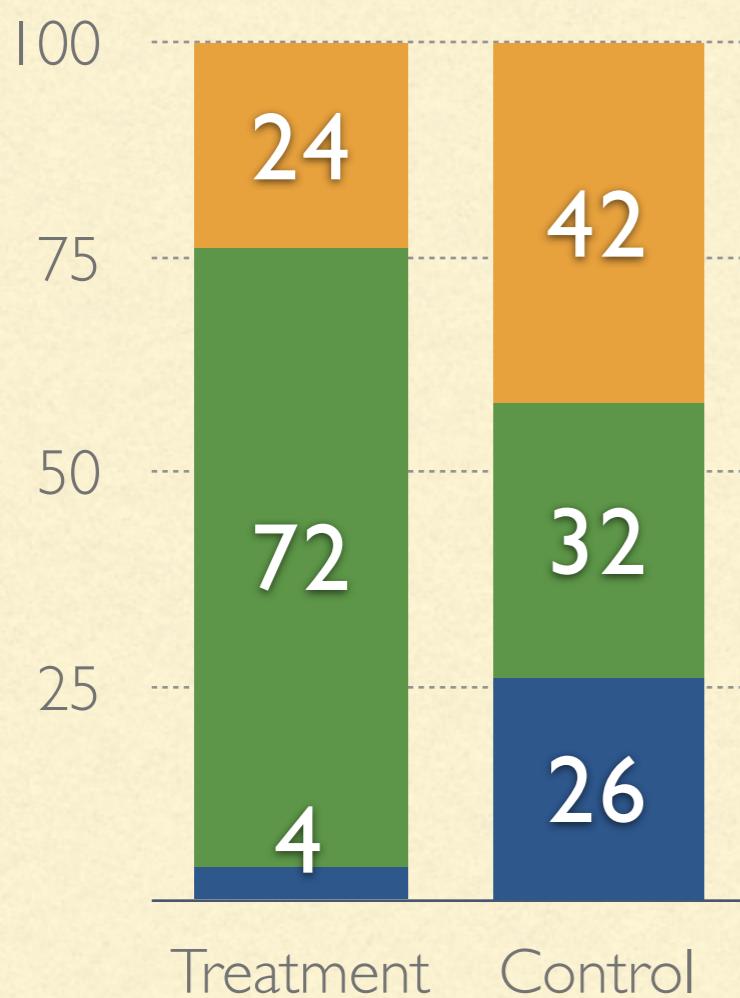
Efficiencies



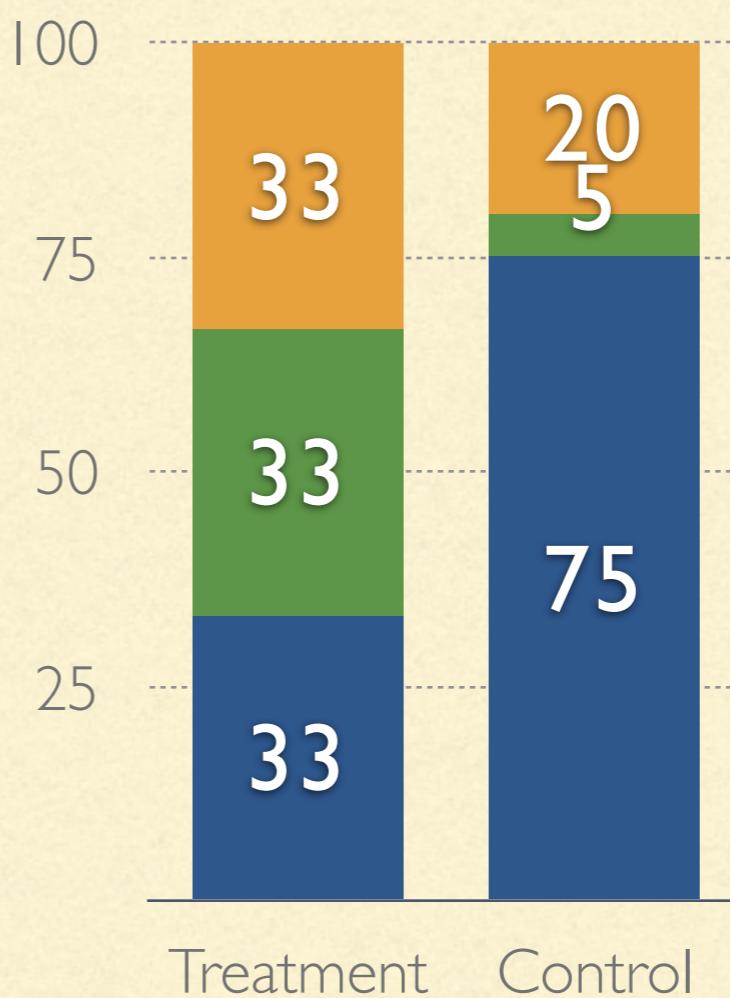
We cannot accurately estimate changes in proportions.

- **Taxon green** is high efficiency; its abundance increased (Ctrl vs Tmt). Additionally, **taxon blue** is low efficiency, and its abundance decreased.
- **Taxon orange**'s abundance depends on the abundance of the other taxa

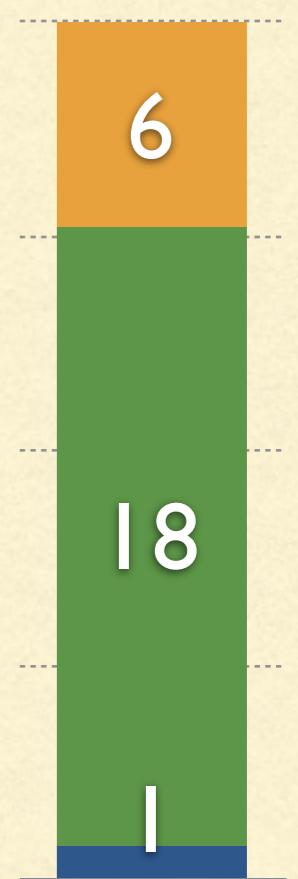
## Observed



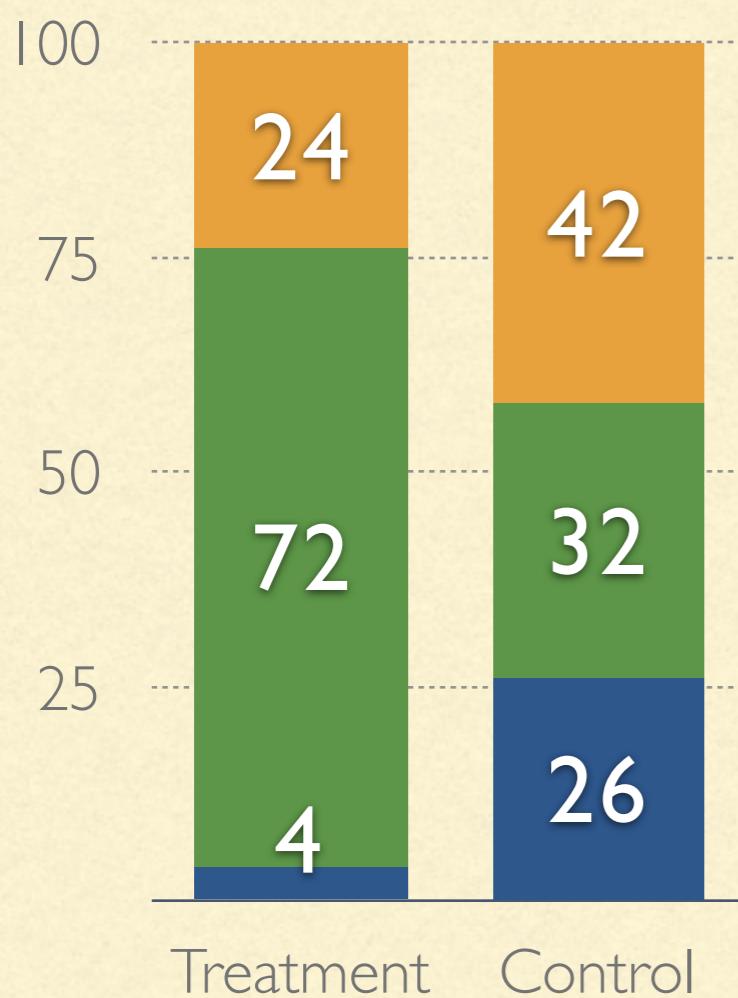
## Actual



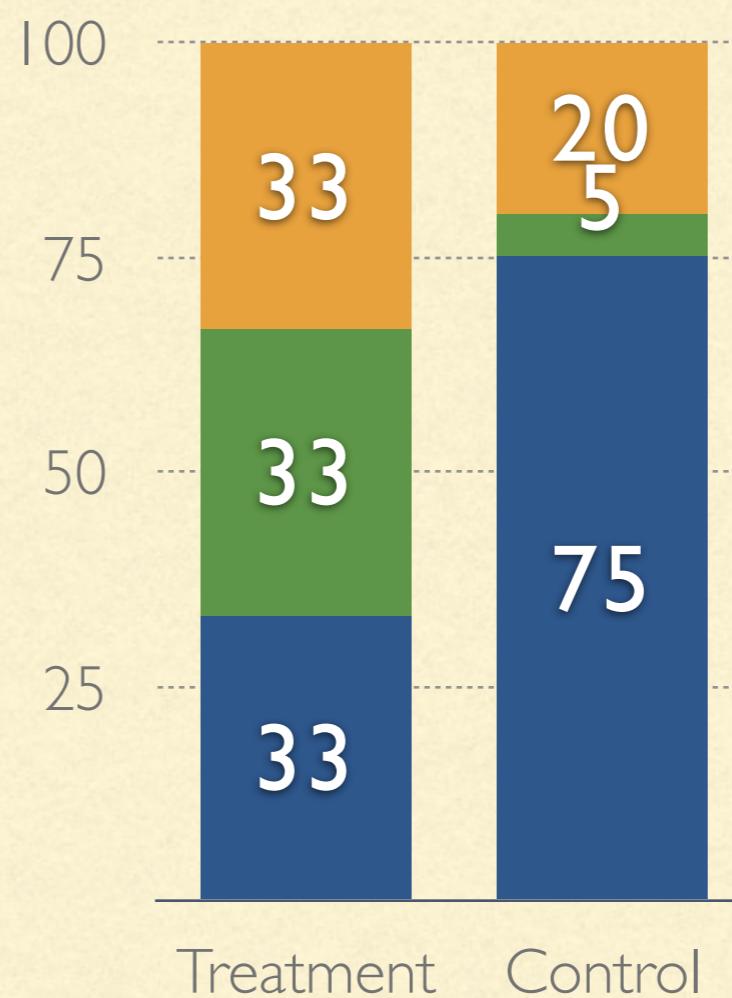
## Efficiencies



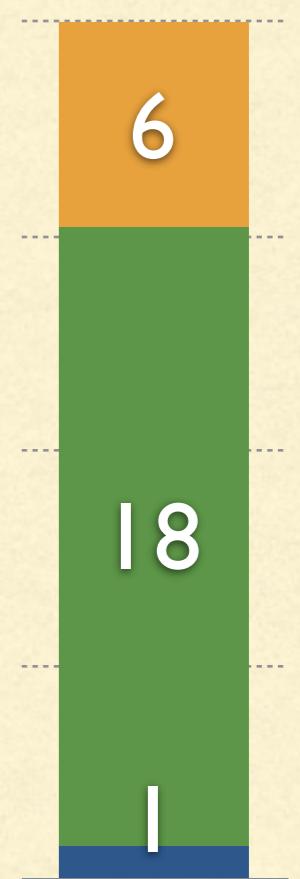
## Observed



## Actual



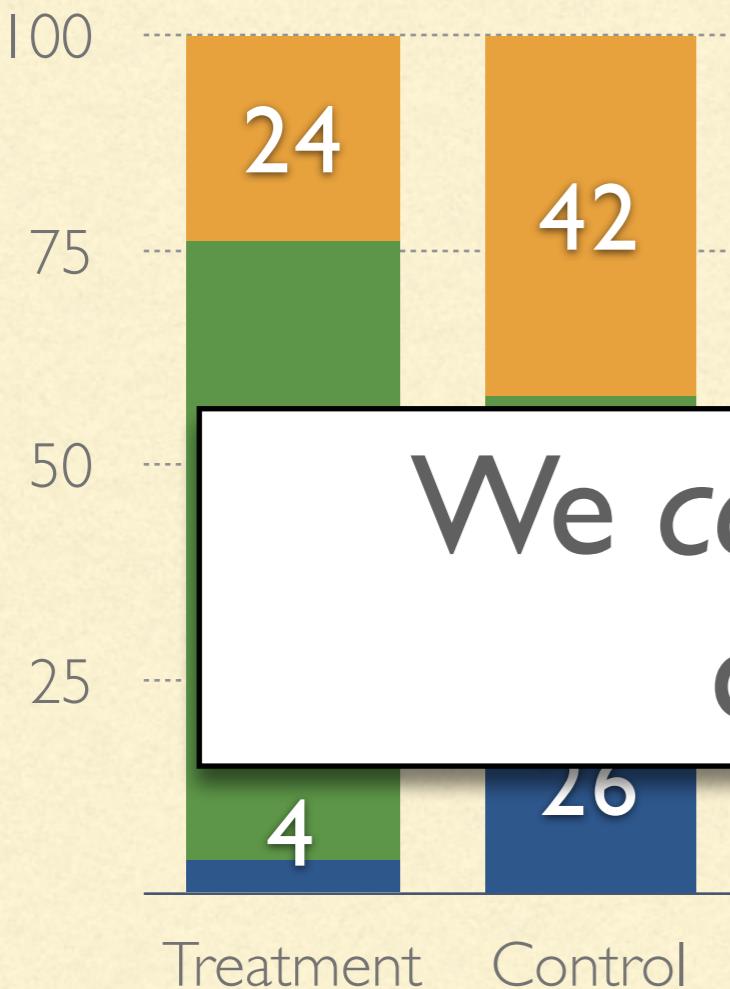
## Efficiencies



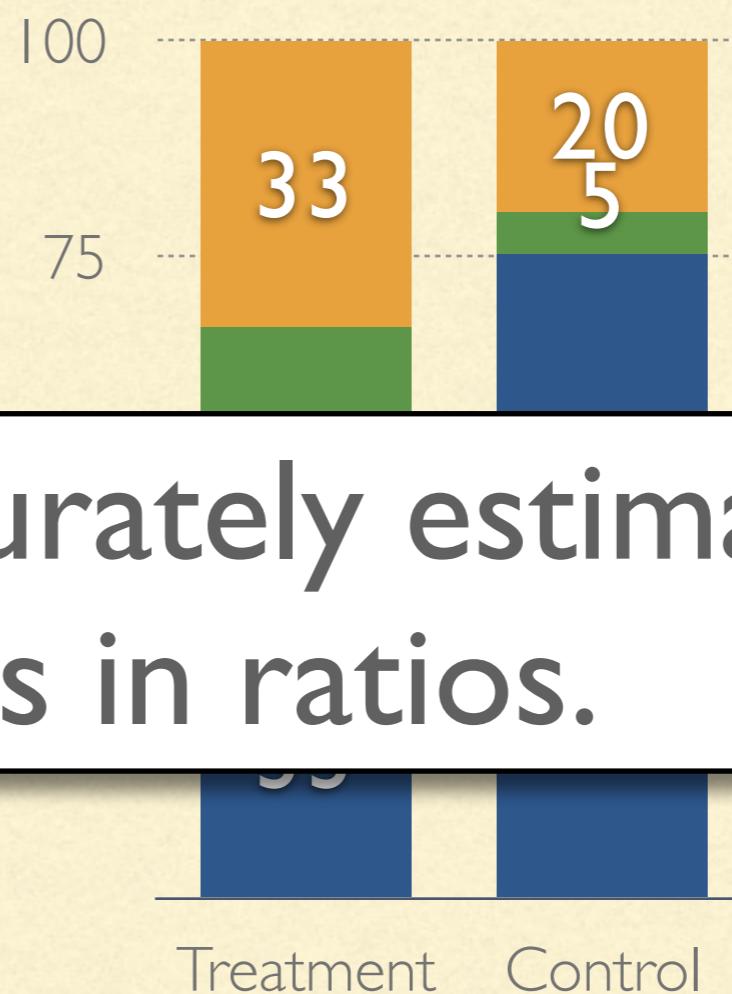
$$\frac{72}{4} \Bigg/ \frac{32}{26} = 15$$

$$\frac{33}{33} \Bigg/ \frac{5}{75} = 15$$

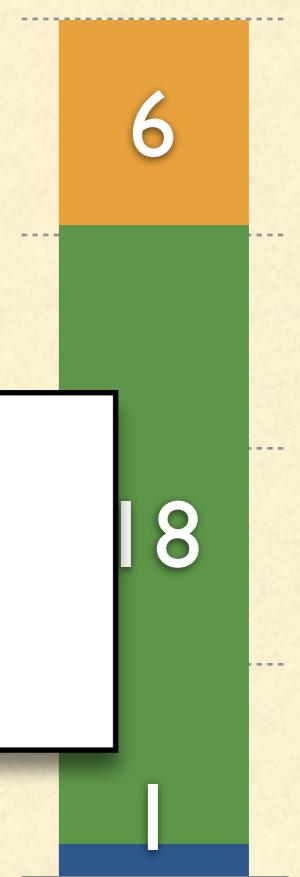
Observed



Actual



Efficiencies



We can accurately estimate changes in ratios.

$$\frac{72}{4} \Bigg/ \frac{32}{26} = 15$$

$$\frac{33}{33} \Bigg/ \frac{5}{75} = 15$$



Michael  
McLaren  
(NCSU)

## 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>



Ben  
Callahan  
(NCSU)

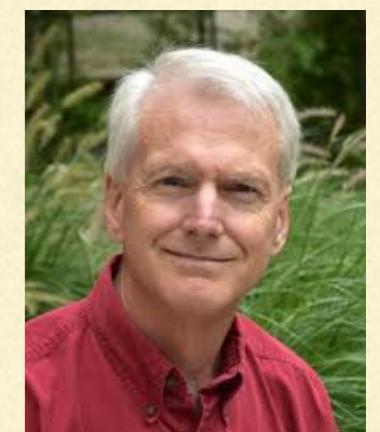
- Proposed model: taxon-specific efficiencies (“biases”) confound estimation of relative abundance
- Model is validated in 16S & shotgun metagenomic data
  - Shotgun data is *also biased* — see paper for details

The image shows a screenshot of the bioRxiv preprint server website. At the top left is the CSHL logo with the text "Cold Spring Harbor Laboratory". The bioRxiv logo is prominently displayed in the center. Below it, the text "THE PREPRINT SERVER FOR BIOLOGY" is visible. On the right side of the header are links for "HOME" and "ABOUT", and a search bar. In the main content area, there are two buttons: "New Results" on the left and "Comment on this paper" on the right. The title of the preprint is "A multi-view model for relative and absolute microbial abundances". Below the title, the authors are listed as "Brian D. Williamson, James P. Hughes, Amy D. Willis". A DOI link "doi: <https://doi.org/10.1101/761486>" is also provided.



Brian  
Williamson  
(UW)

- Calibration with taxon-specific qPCR + 16S
  - Combing low-throughput + high-throughput
  - Key implication: Do not multiply total 16S concentration and 16S relative abundance to estimate taxon-specific concentration



Jim Hughes  
(UW)

# IMPLICATIONS

---

- Analyse ratios, not proportions
  - If you must analyse proportions, focus on largest signals
- Avoid aggregating abundances by taxonomy/phylogeny
- Avoid “quantitative” diversity indices

# ONGOING & FUTURE WORK

---

- Robust statistical methods
- Protocol selection
- Optimal experimental design (positive + negative controls)
- Bias prediction from phylogeny
- Calibration & meta-analysis (?)

# SUMMARY

---

- Proposed a mathematical description of bias in sequencing
  - Validated on 16S and shotgun metagenomics data
- Model has serious implications for analysing microbiome data
  - *Proportions can lead to conclusions in the wrong direction*
  - Certain analysis methods are robust to bias; others are not

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>



New Results

Comment on this paper

**A multi-view model for relative and absolute microbial abundances**

Brian D. Williamson, James P. Hughes, Amy D. Willis  
**doi:** <https://doi.org/10.1101/761486>



Michael  
McLaren  
(NCSU)



Ben  
Callahan  
(NCSU)



David  
Clausen  
(UW)



Jim  
Hughes  
(UW)



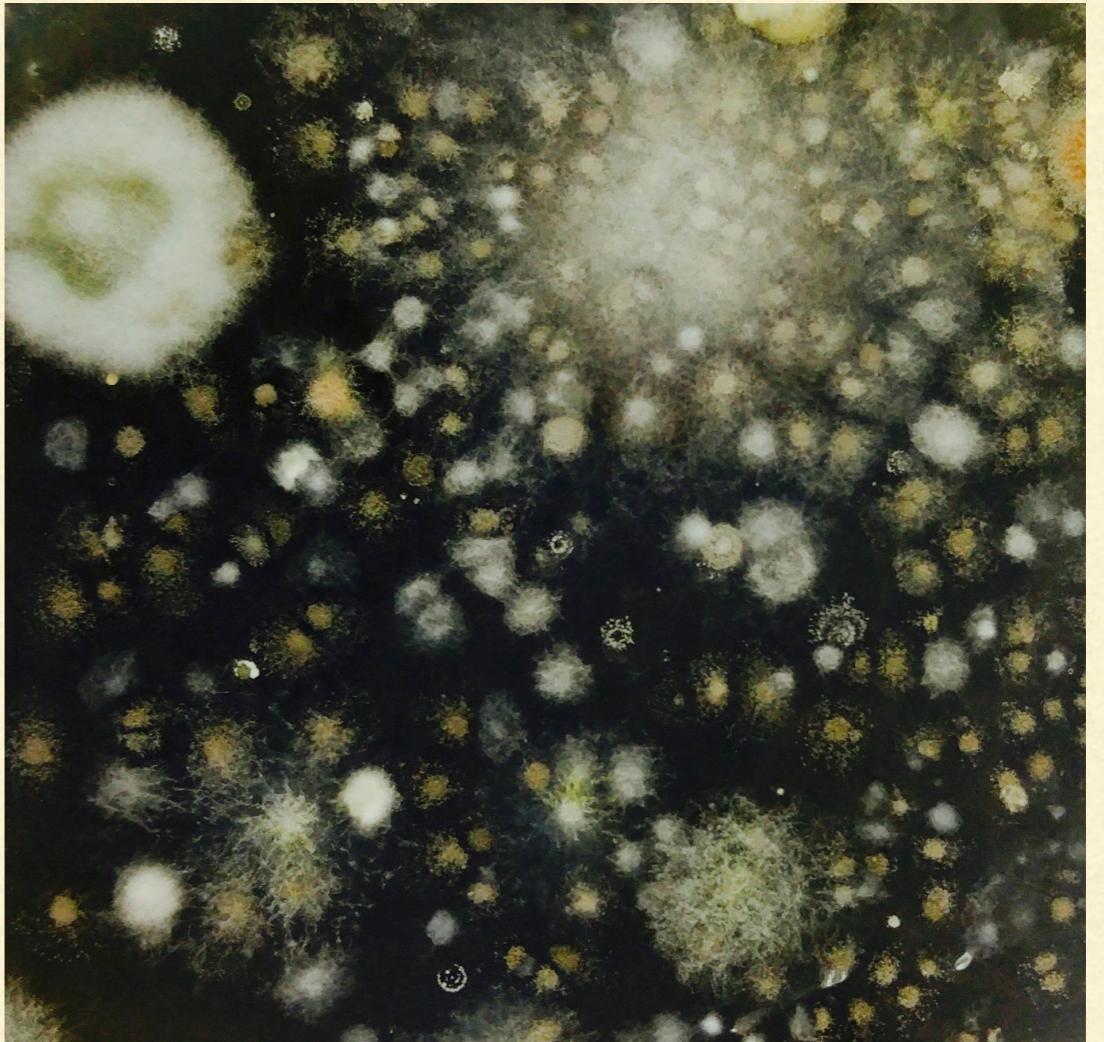
Brian  
Williamson  
(UW)<sup>28</sup>

# MEASUREMENT ERROR IN MICROBIOME DATA

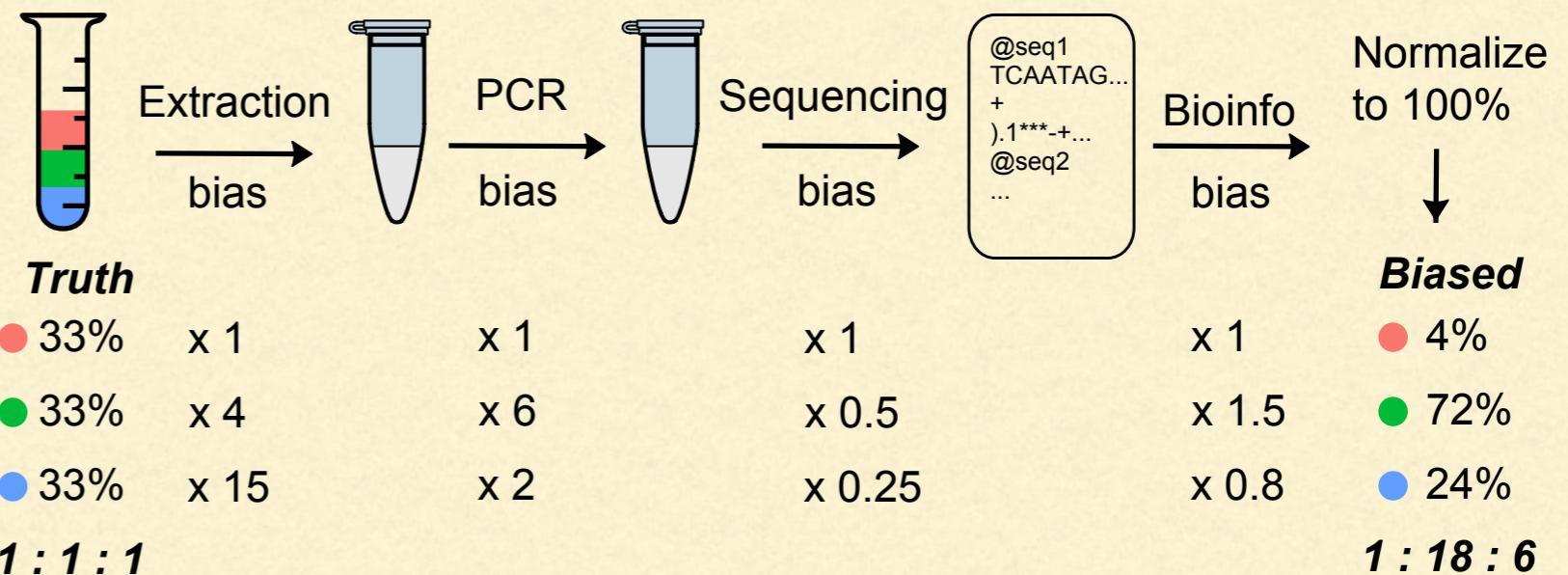
Amy D Willis PhD  
Assistant Professor  
Department of Biostatistics  
University of Washington, USA



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



# SUPPLEMENTARY SLIDES



Amy D Willis PhD  
Assistant Professor  
Department of Biostatistics  
University of Washington, USA

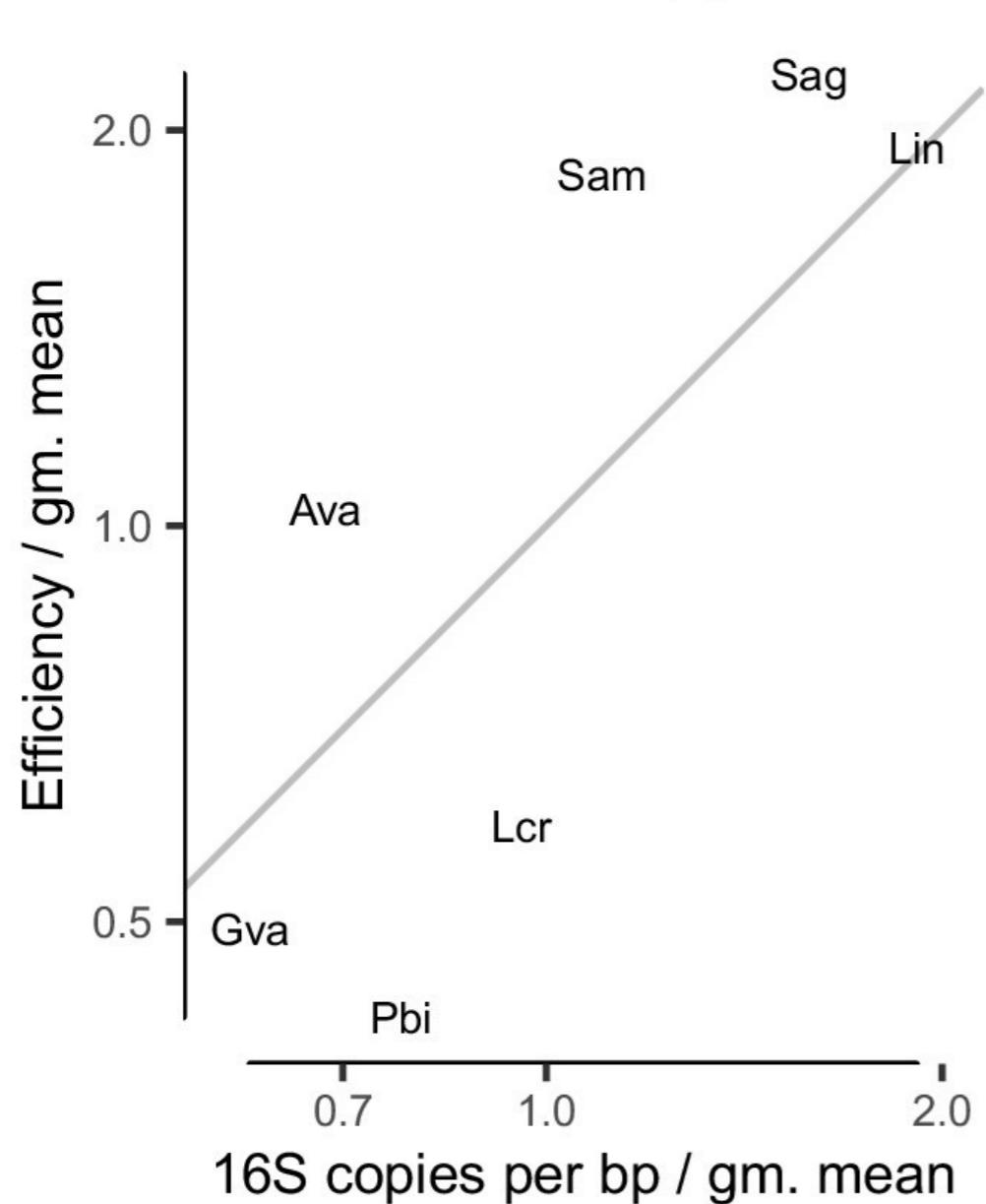
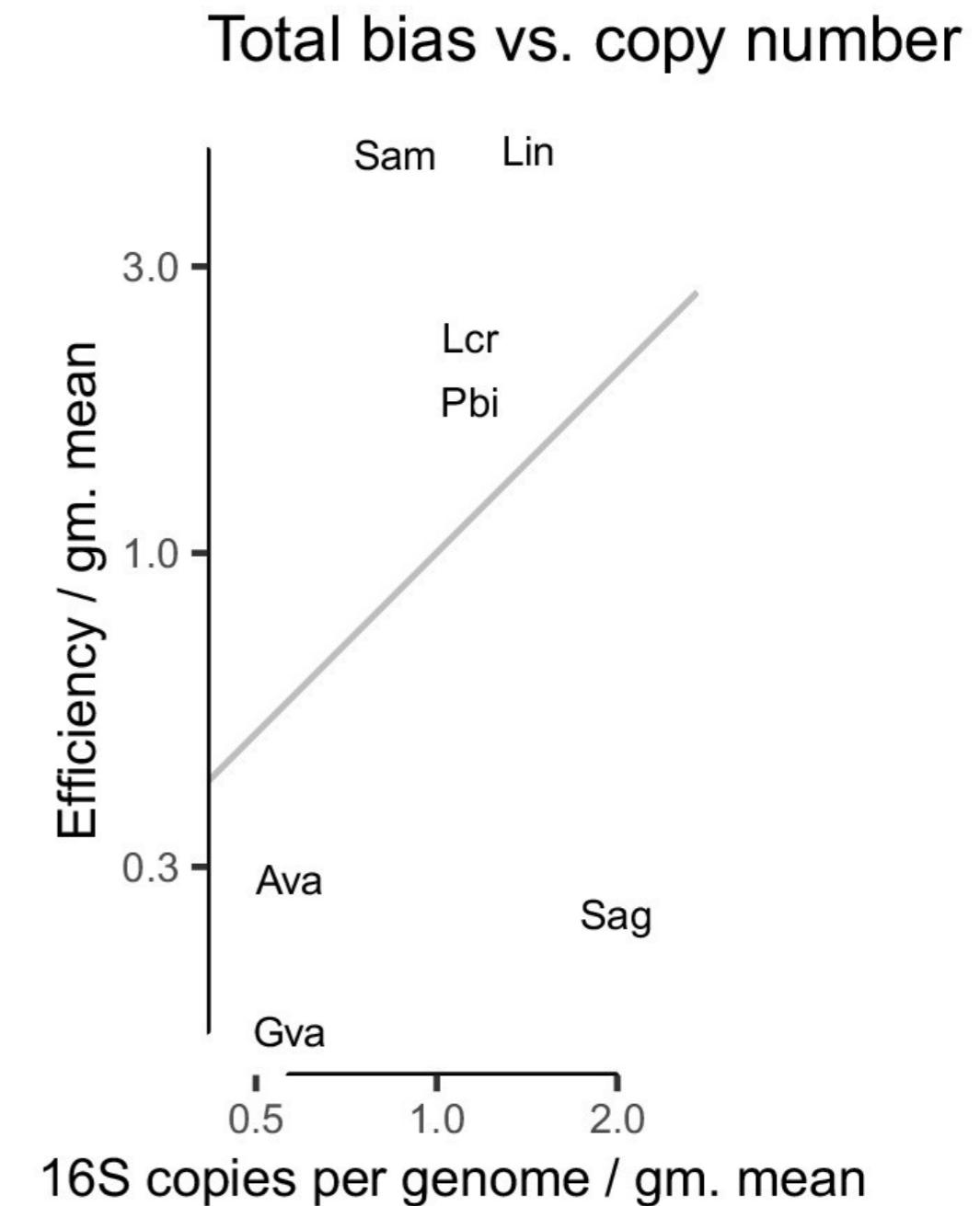


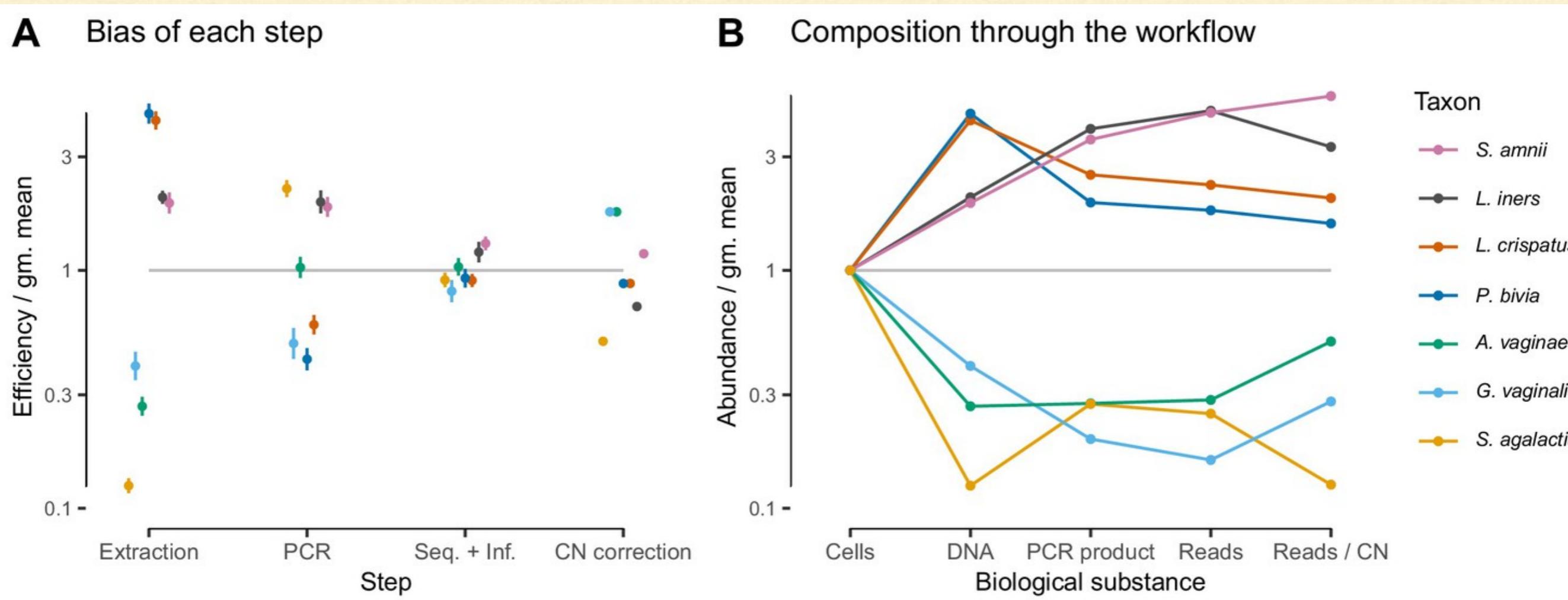
@AmyDWillis



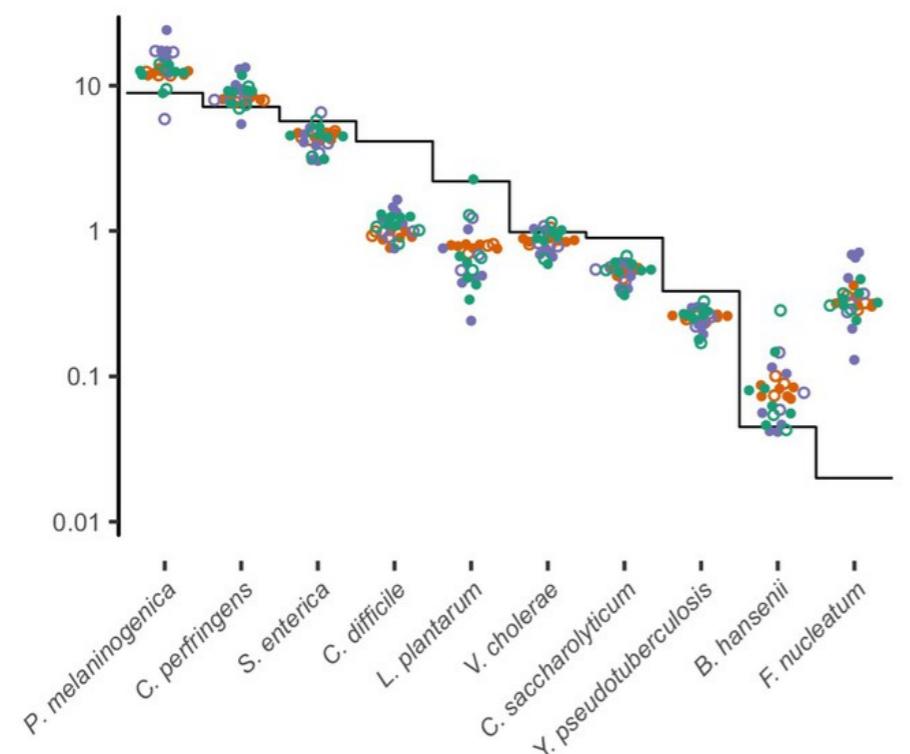
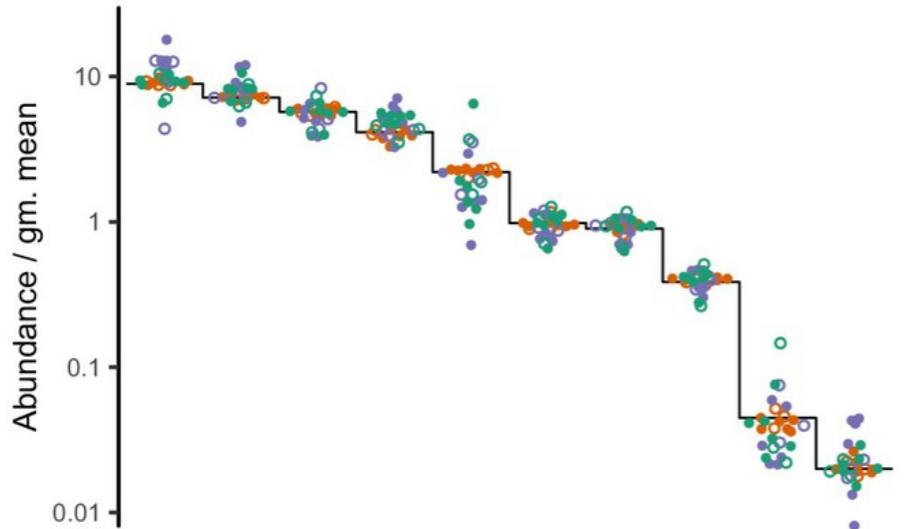
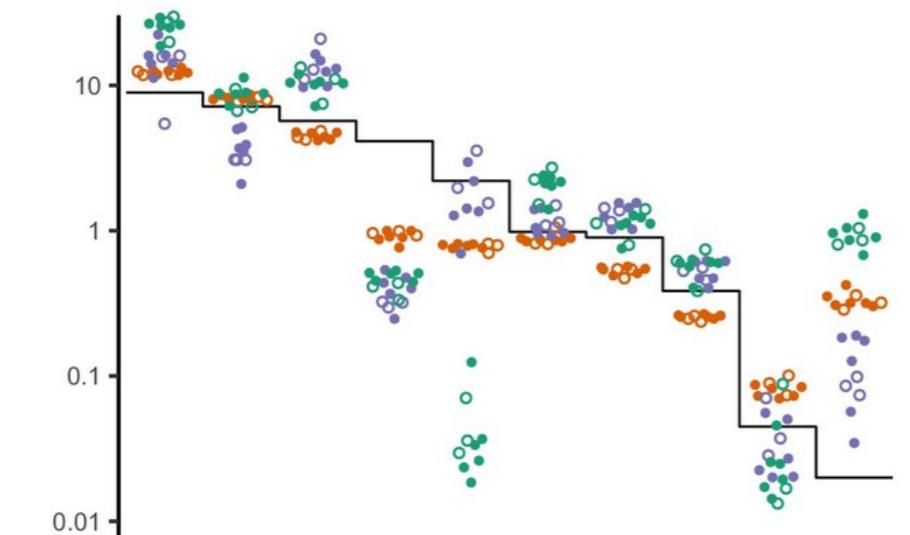
[adwillis@uw.edu](mailto:adwillis@uw.edu)



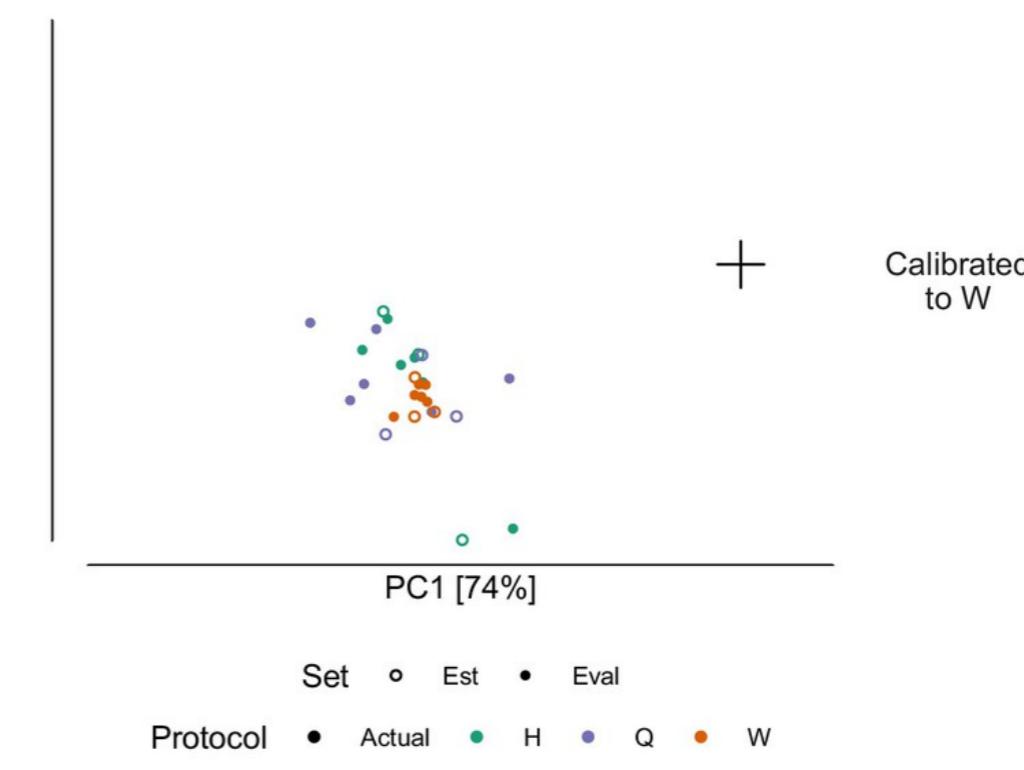
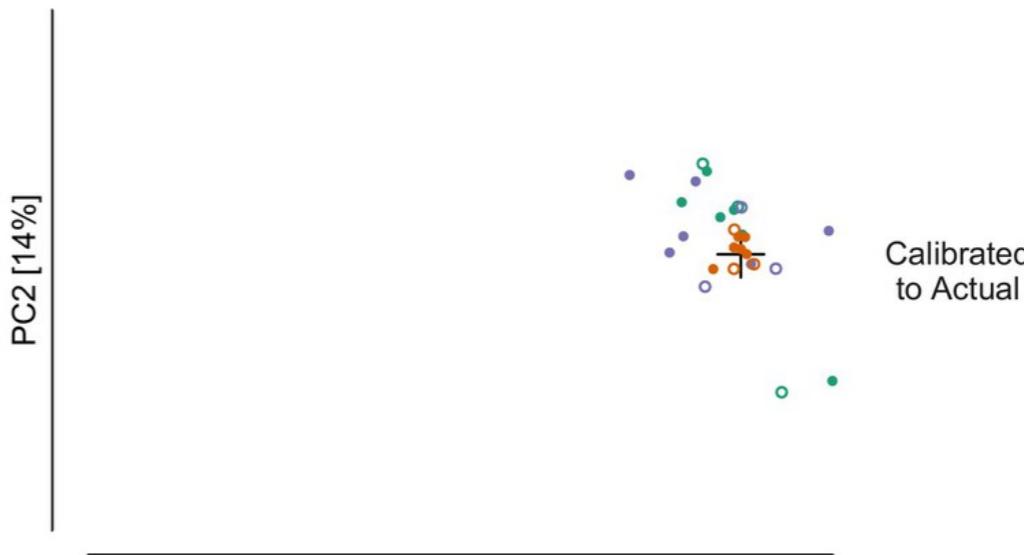
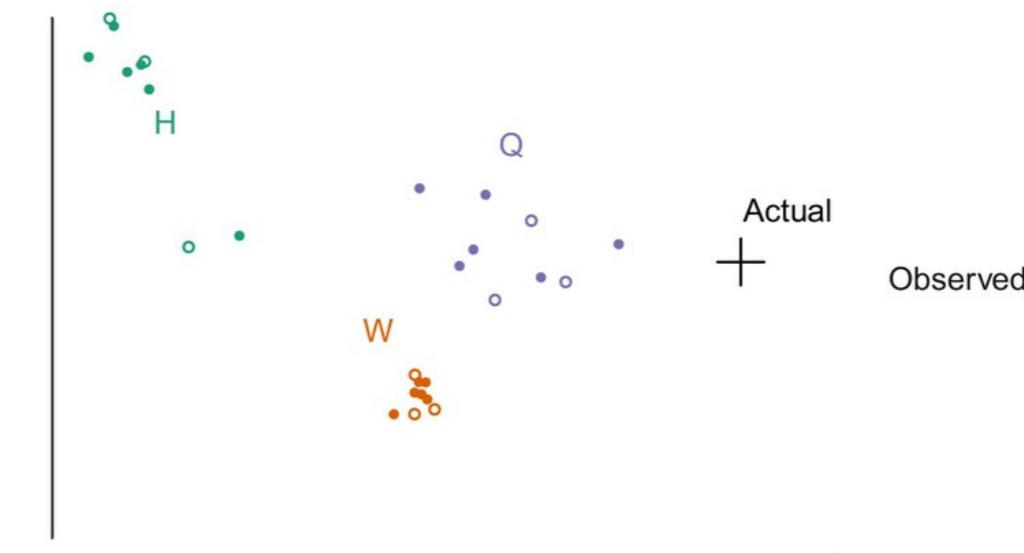
**A****B**

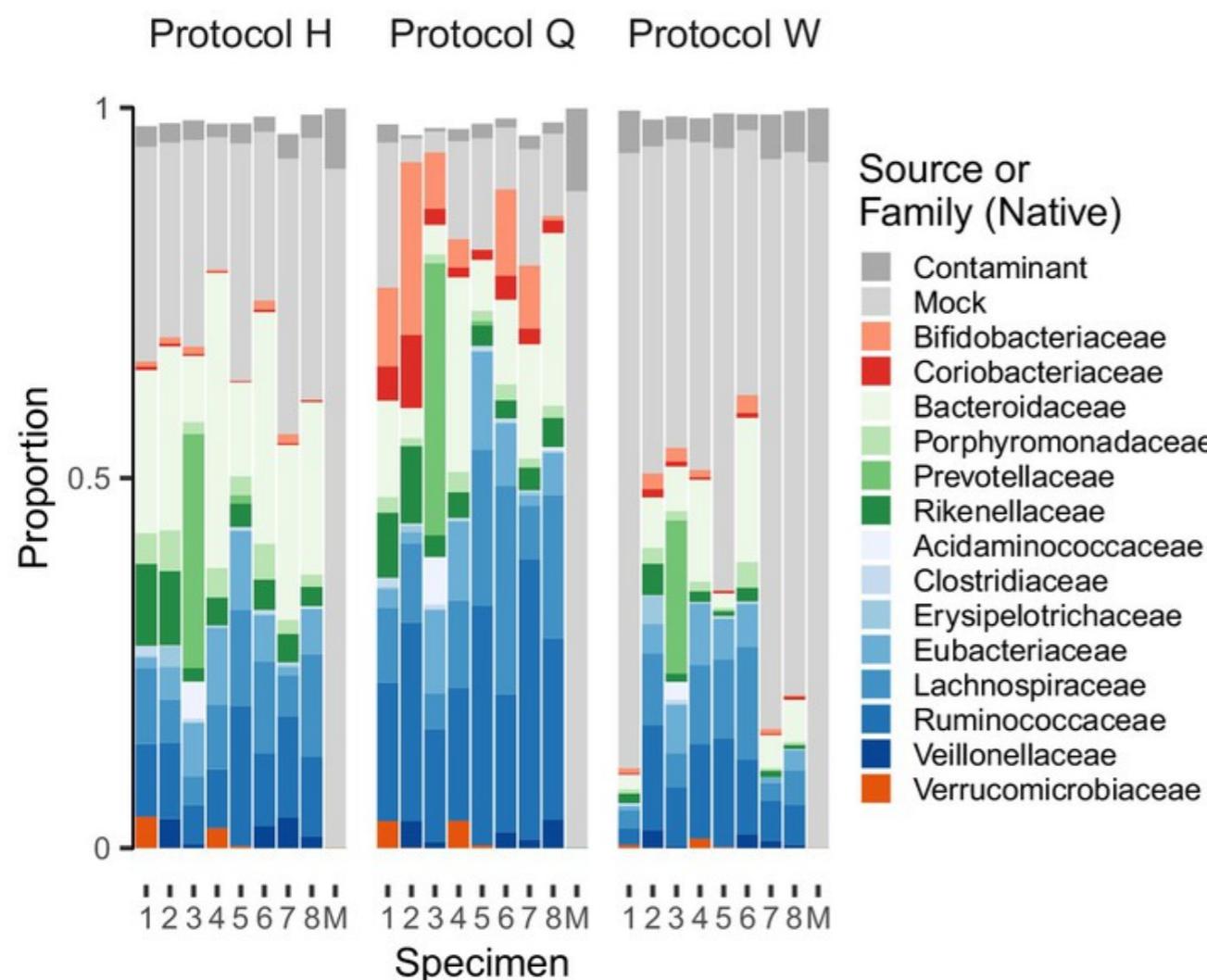
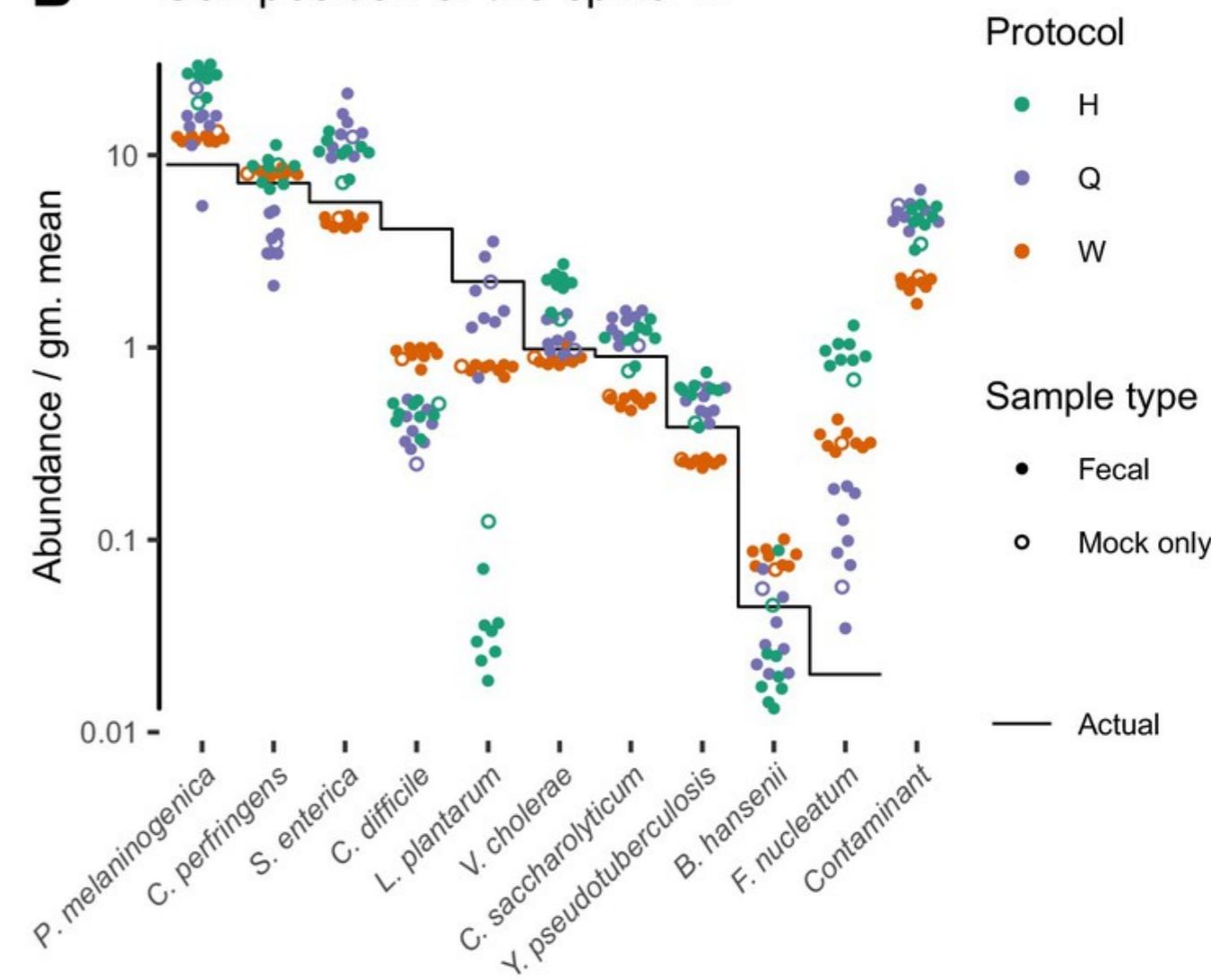


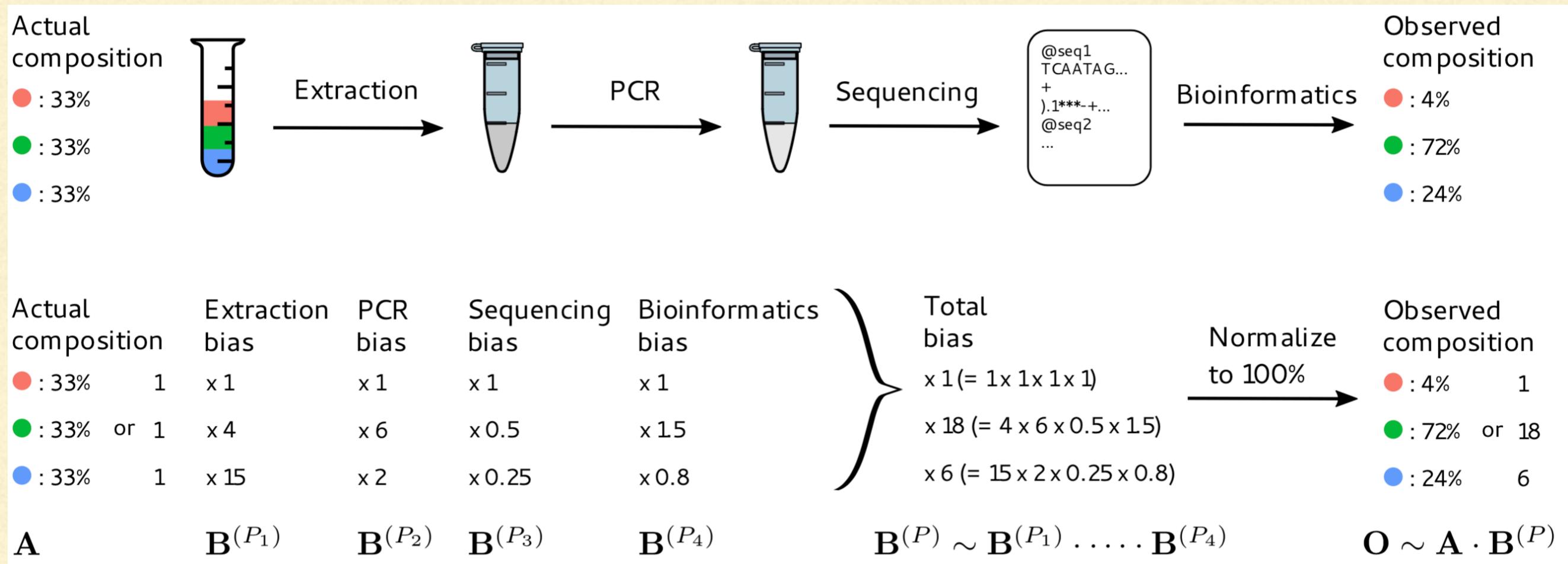
Taxon relative abundances

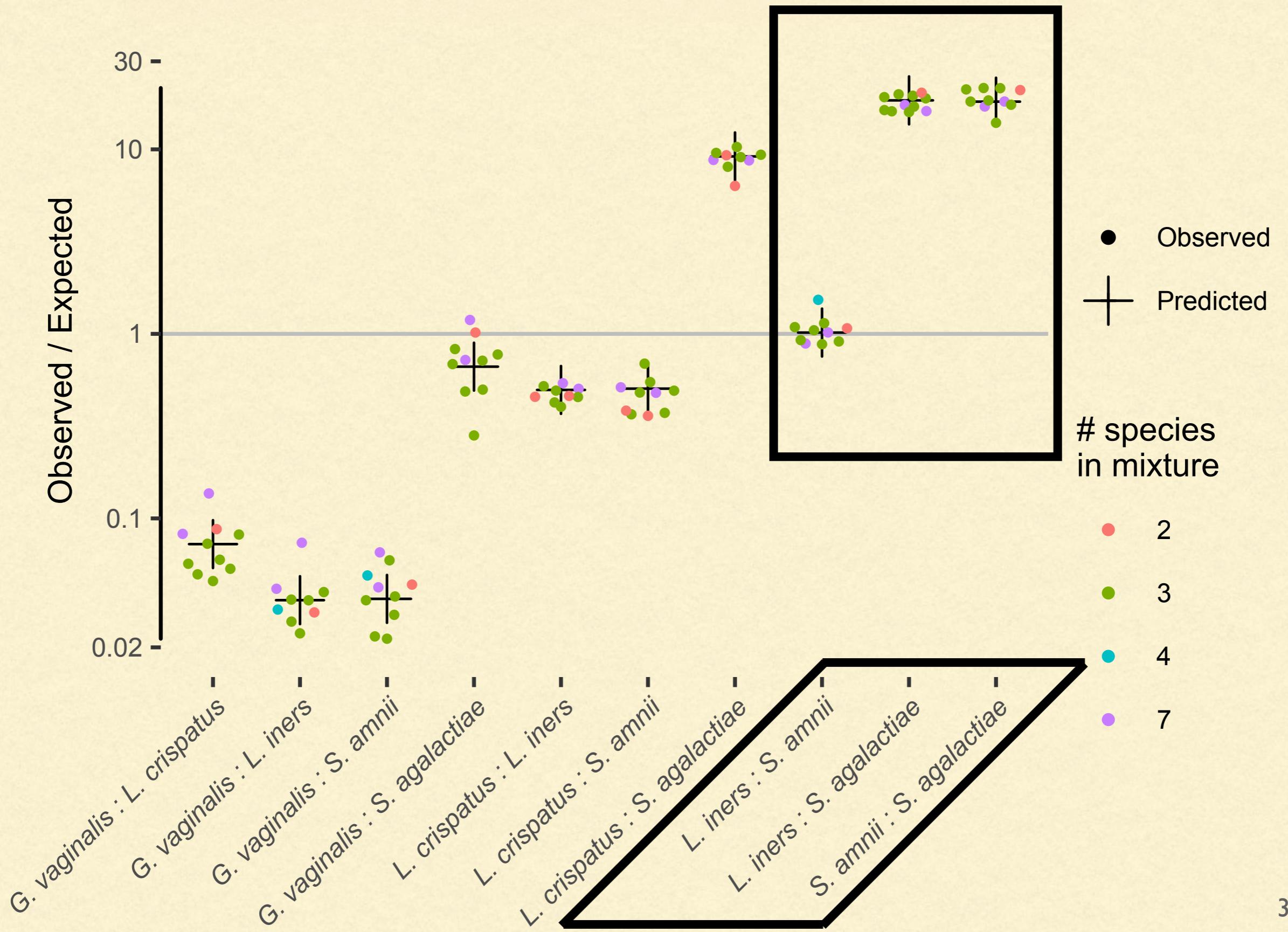


Sample ordination



**A Composition of all bacteria****B Composition of the spike-in**





# WHY ARE RATIOS OF RATIOS UNAFFECTED BY BIAS?

---

$$\left( \frac{\frac{p_{11}e_1}{\sum p_{1j}e_j}}{\frac{p_{12}e_2}{\sum p_{1j}e_j}} \right) \left/ \left( \frac{\frac{p_{21}e_1}{\sum p_{2j}e_j}}{\frac{p_{22}e_2}{\sum p_{2j}e_j}} \right) \right. = \frac{p_{11}e_1}{p_{12}e_2} \left/ \left( \frac{p_{21}e_1}{p_{22}e_2} \right) \right. = \frac{p_{11}}{p_{12}} \left/ \left( \frac{p_{21}}{p_{22}} \right) \right.$$

$p_{ij}$  = proportion of taxon j in sample i  
 $e_j$  = efficiency of taxon j

# HOW CONSISTENT IS MICROBIOME DATA WITHIN AND ACROSS LABS?

