

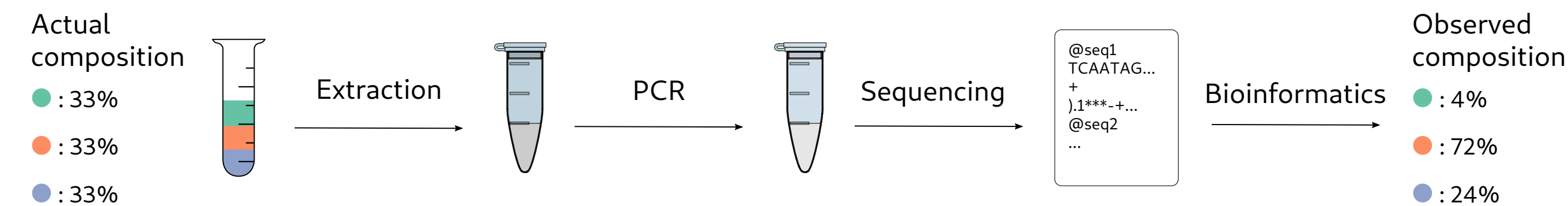
Consistent and correctable bias in metagenomic sequencing measurements

Michael R. McLaren¹, Amy D. Willis², and Benjamin J. Callahan¹

¹NC State University ²University of Washington

A mathematical model of bias

Each step in a marker-gene or metagenomics workflow is **biased** towards detecting certain taxa over others



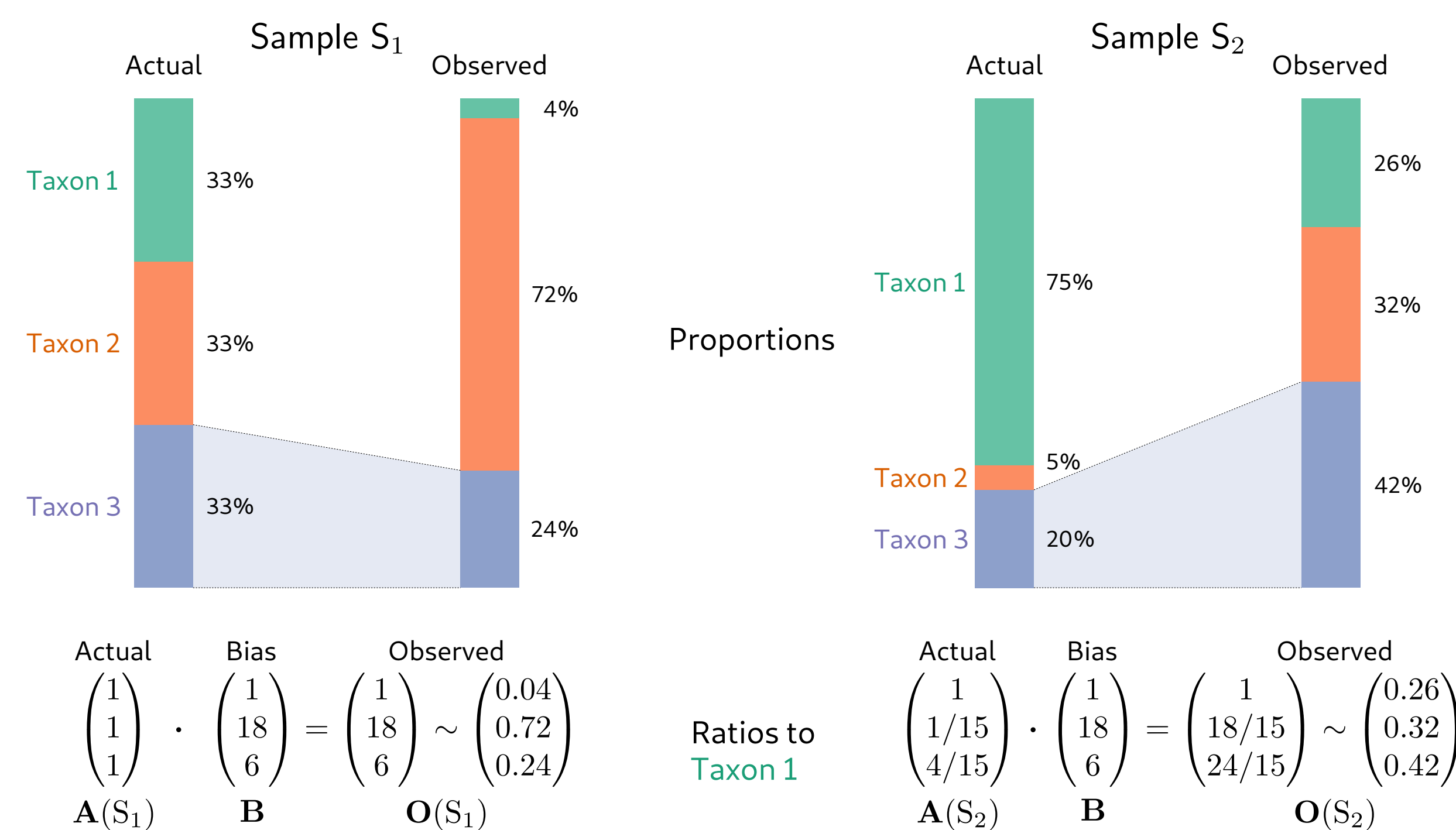
Bias makes measured abundances inaccurate and quantitatively incomparable between experiments. The input-output (actual-observed) relationship is currently unknown

Model of bias: Detection efficiencies multiply across steps

Actual composition	Extraction bias	PCR bias	Sequencing bias	Bioinformatics bias	Total bias	Normalize to 100%	Observed composition
● : 33%	1	x1	x1	x1	x1 (= 1 x 1 x 1 x 1)	● : 4%	1
● : 33% or 1	x4	x6	x0.5	x1.5	x18 (= 4 x 6 x 0.5 x 1.5)	● : 72% or 18	1
● : 33%	1	x15	x2	x0.8	x6 (= 15 x 2 x 0.25 x 0.8)	● : 24%	6

The set of **relative detection efficiencies** for all taxa for a given protocol determines the bias of that protocol for arbitrarily composed samples

Bias acts consistently on taxon ratios, not taxon proportions



Consistent bias can be estimated and corrected (**calibration**)

Estimation from a single control sample

$$\hat{B} \sim O(c)/A(c)$$

Estimation from multiple controls containing different, partially-overlapping sets of taxa

$$\hat{B} = \arg \min_B \sum_c \|O(c)/(A(c) \cdot B)\|^2$$

Calibration to actual composition

$$\hat{A}(s) \sim O(s)/\hat{B}$$

Calibration to a reference protocol R

$$\hat{B}^{(P/R)} \sim O^{(P)}(c)/O^{(R)}(c)$$

$$\hat{O}^{(P)}(s) \sim O^{(R)}(s)/\hat{B}^{(P/R)}$$

Evaluation datasets

16S amplicon data from Brooks et al (2015)

7 bacterial species in 58 unique mixtures of cells, DNA, and PCR product test that bias is independent of which species are in the sample

Mixtures of cells, of DNA, and of PCR product test the model for extraction, PCR, and the total workflow

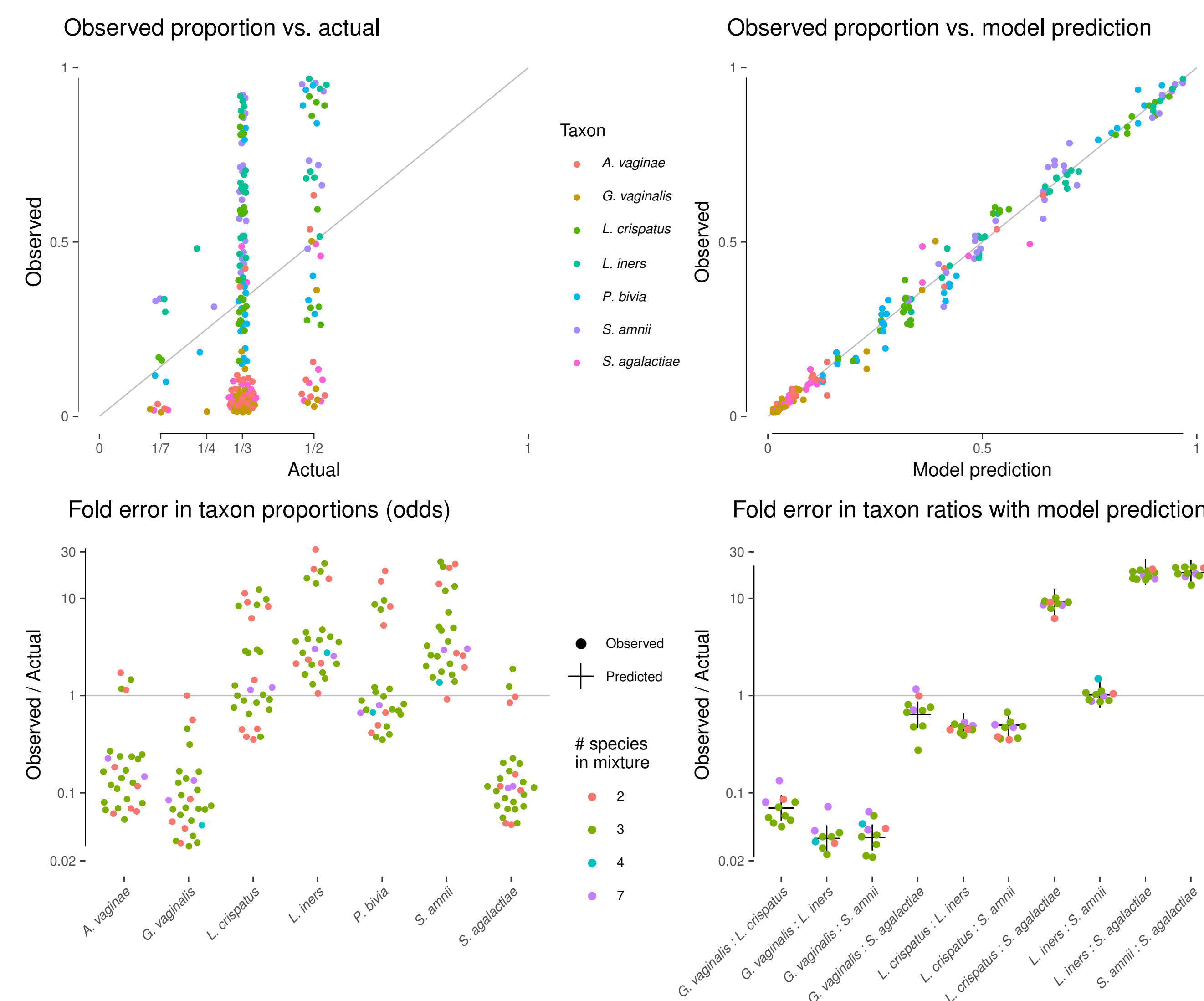
Shotgun data from Costea et al (2017)

10 bacterial species in fixed amounts measured alone and spiked into feces test that bias of spike-in is independent of the background composition

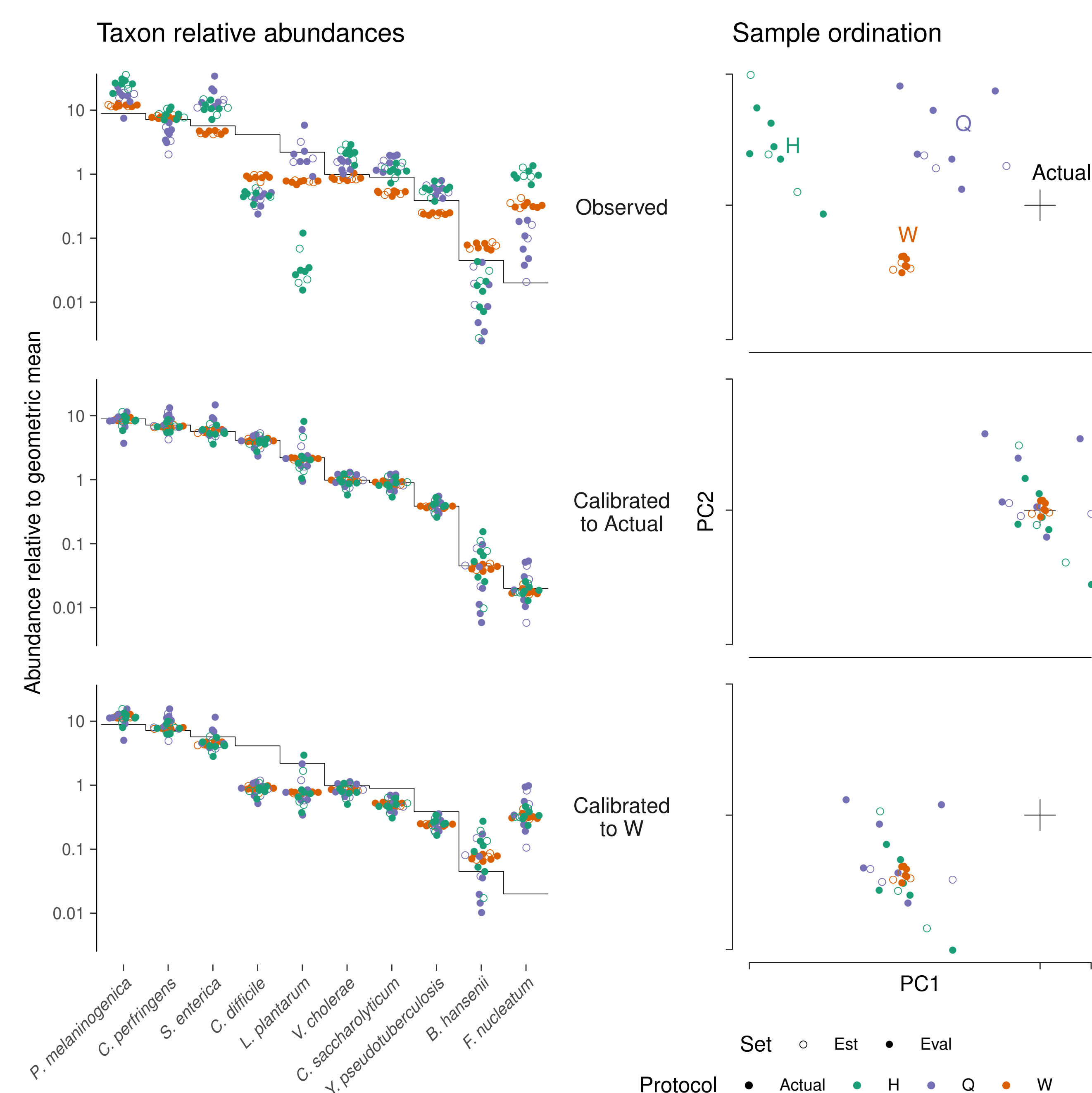
3 DNA extraction protocols (H, Q, and W) compare calibration to actual composition and calibration to a reference protocol

Results

Model performance in the 16S-sequenced cell mixtures



Bias and calibration in the shotgun-sequenced spike-ins

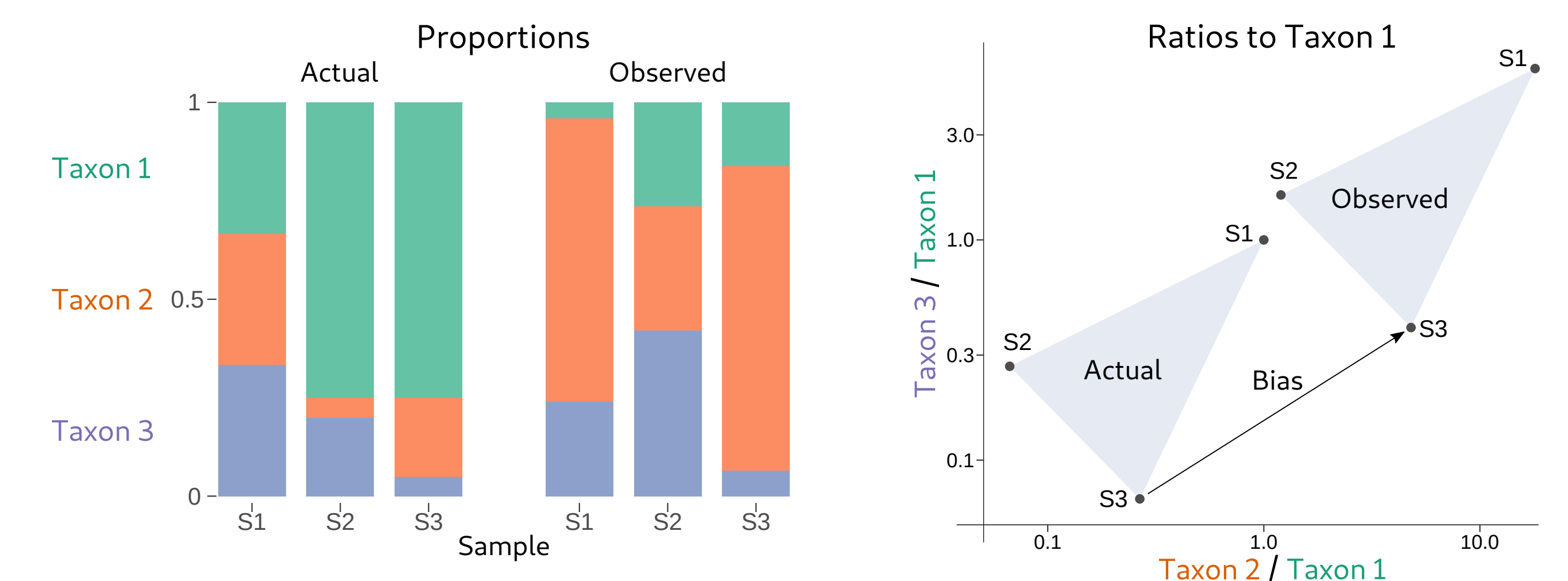


Next steps

Which (meta-)analyses are more robust to bias?

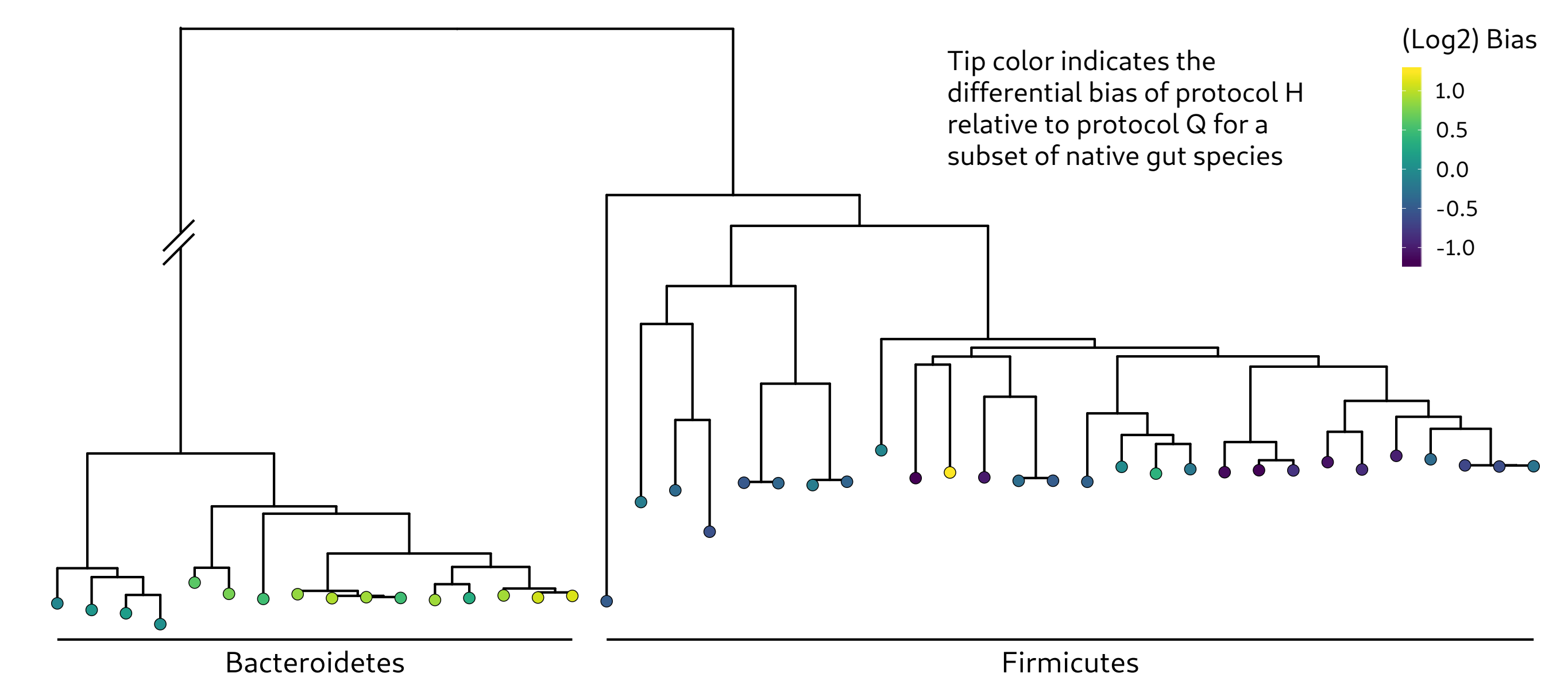
Analyses based on taxon ratios (e.g., CoDA methods) are theoretically bias invariant

$$\frac{O_i(S1)/O_j(S1)}{O_i(S2)/O_j(S2)} = \frac{A_i(S1)B_j}{A_j(S1)B_i} / \frac{A_i(S2)B_j}{A_j(S2)B_i} = \frac{A_i(S1)}{A_j(S1)} / \frac{A_i(S2)}{A_j(S2)}$$



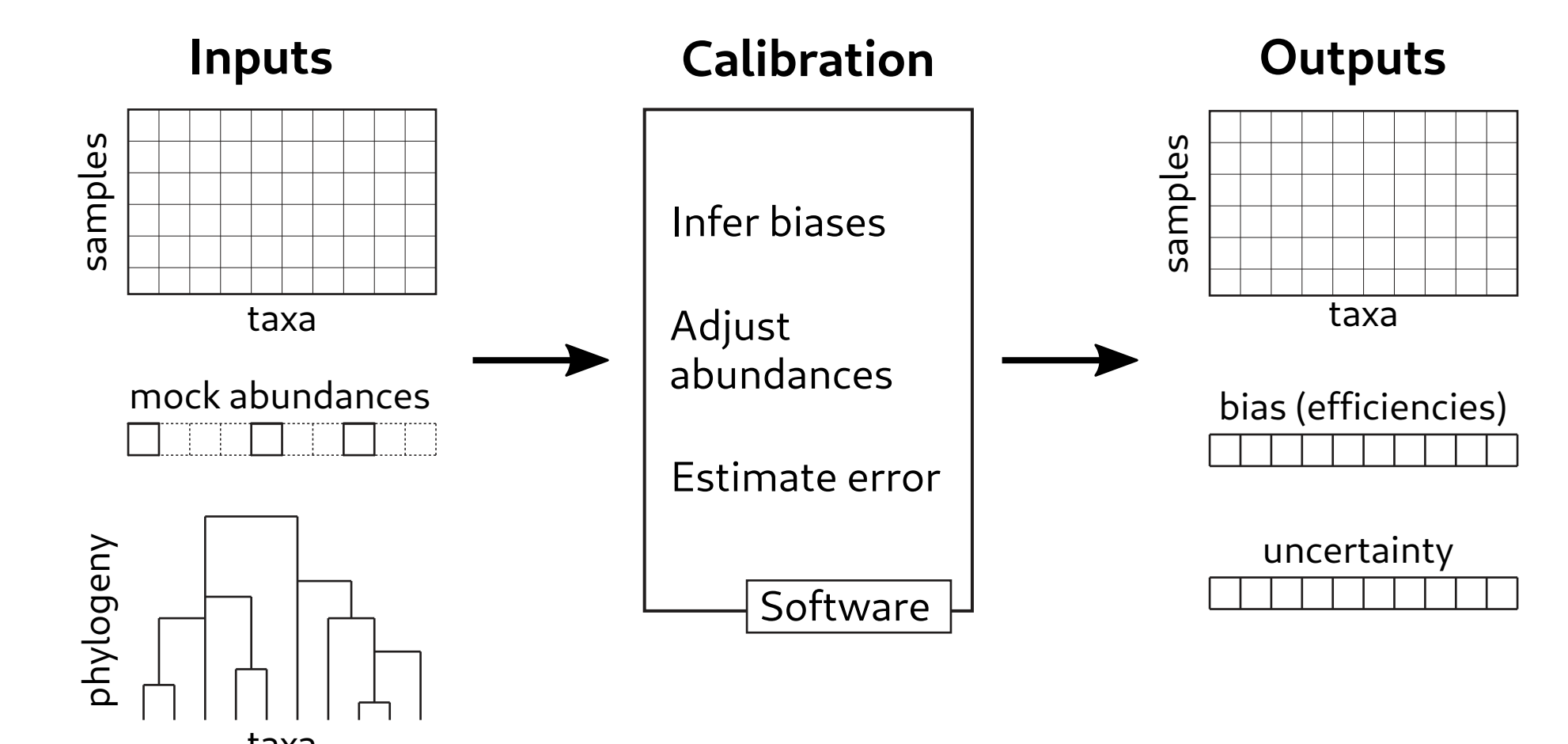
But factors such as biological zeros, cross-sample contamination, and taxonomic aggregation complicate this picture

How does bias evolve?



Natural communities often contain hundreds of species. If bias evolves slowly, we need only measure bias for a small fraction

Long-term goal: Software to estimate the bias of all taxa and calibrate the measurements of complex communities



References

- McLaren MR, Willis AD, Callahan BJ. 2019. Consistent and correctable bias in metagenomic sequencing measurements. bioRxiv 559831.
- Brooks JP, et al. 2015. The truth about metagenomics: quantifying and counteracting bias in 16S rRNA studies. BMC Microbiol 15:66.
- Costea PI, et al. 2017. Towards standards for human fecal sample processing in metagenomic studies. Nat Biotechnol 35:1069-1076.