

Extracting the  
**signal** from the ***noise***  
in high-throughput amplicon data

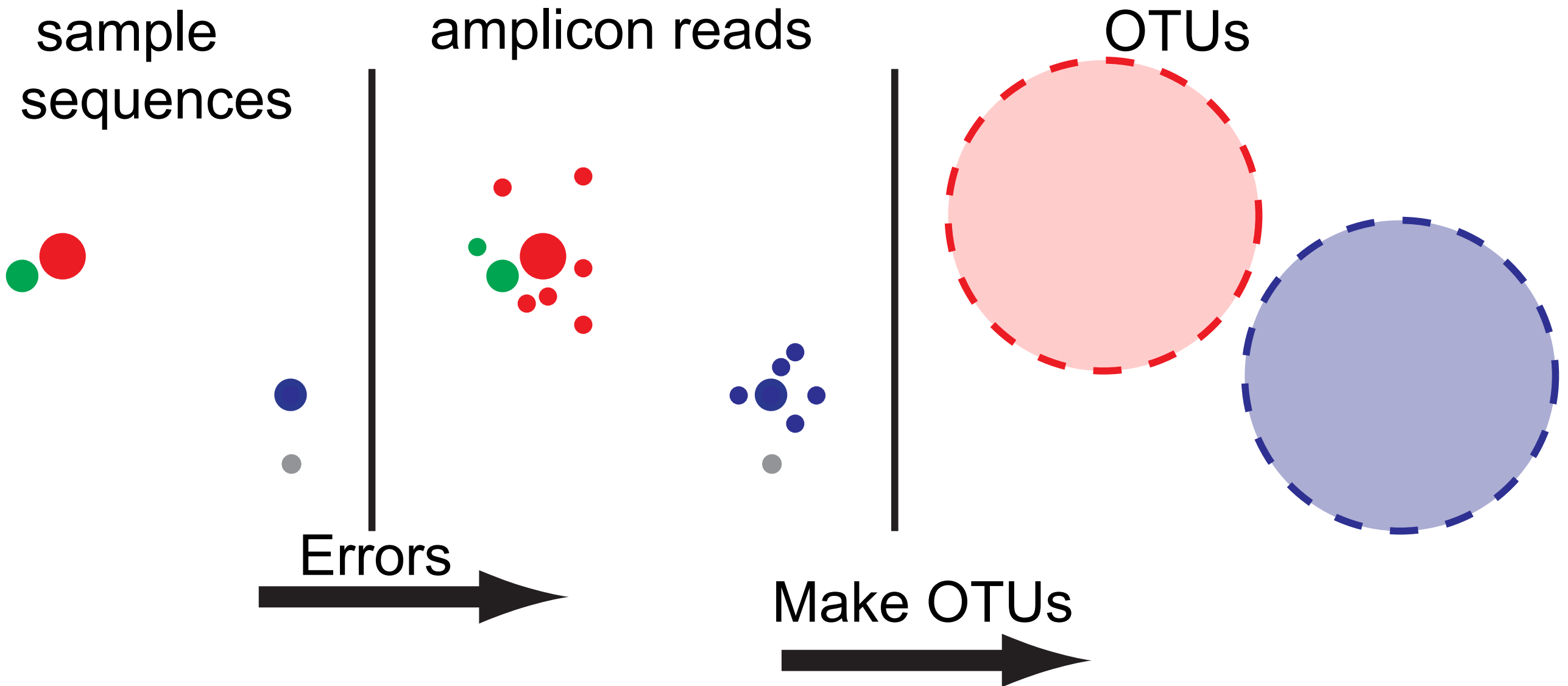
# **The amplicon inference problem**

Infer the sample types and abundances  $\{(s, a)\}$   
from error-ful amplicon reads  $\{r\}$ .

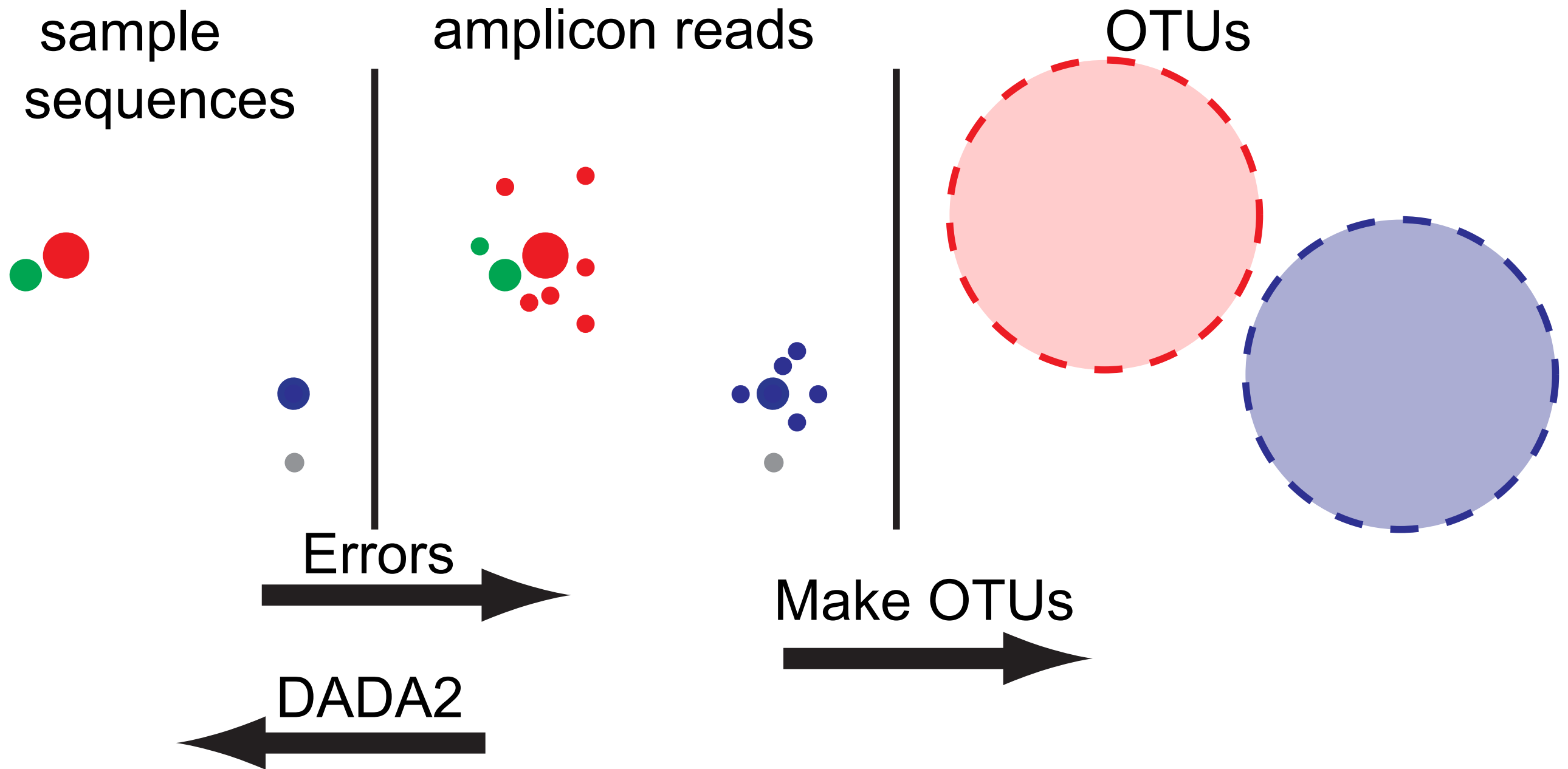
# Amplicon inference problems

- 1. Low resolution**
  - 97%, genus at best
- 2. High false positive rate**
  - #(OTUs)  $\gg$  richness
- 3. Big data scaling**
  - time scales super-linearly
- 4. Cross-study comparison**
  - must reprocess all data together

# DADA2: High resolution



# DADA2: High resolution



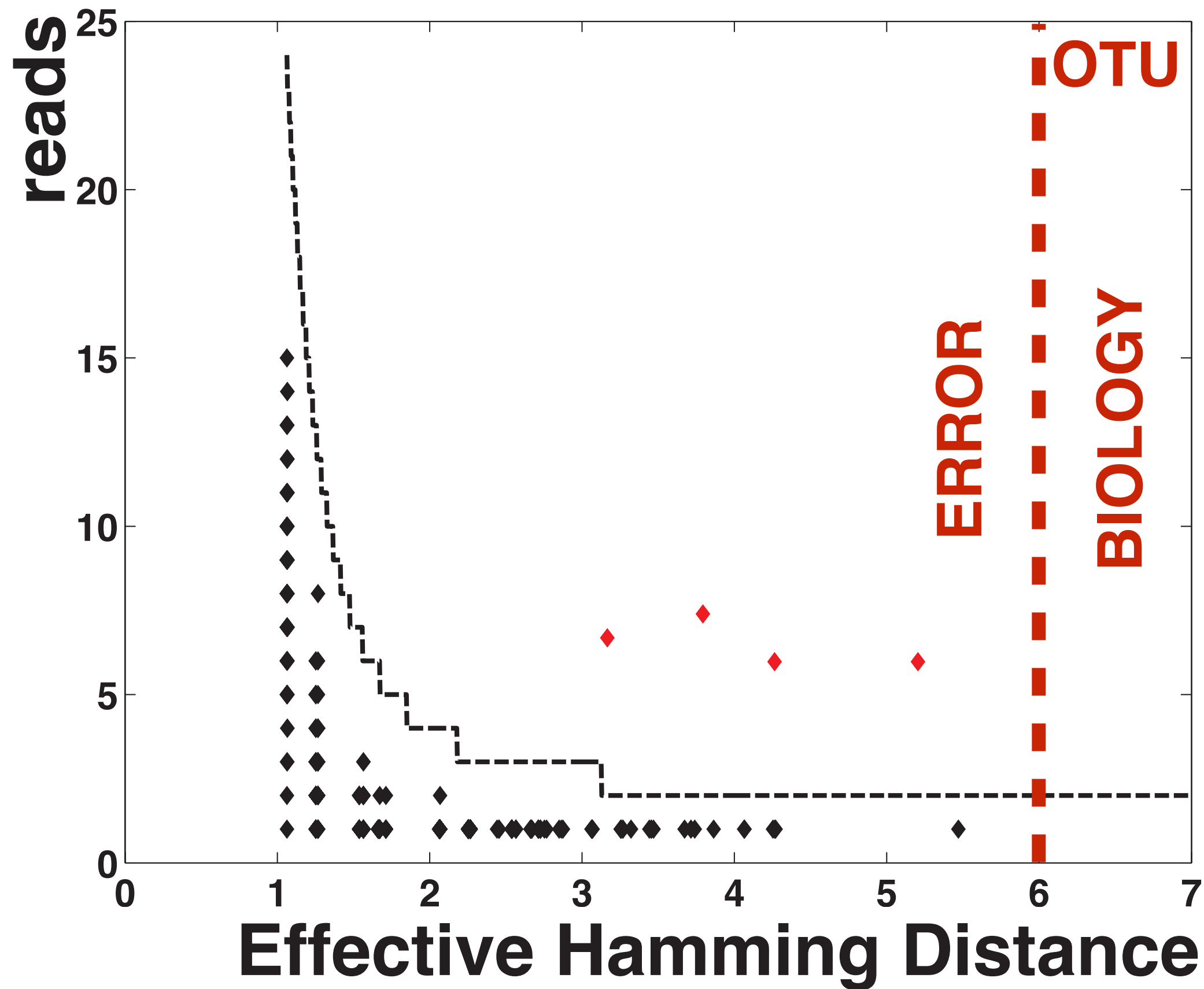
# DADA2: Error model

**s** : ATTAACGAGATTATAACCCAGAGTACGAATA . . .  
      |                                          |  
**r** : ATCAACGAGATTATAACAAGAGTACGAATA . . .

$$p(r|s) = \prod_{i=1}^L p(r(i)|s(i), q_r(i), Z)$$

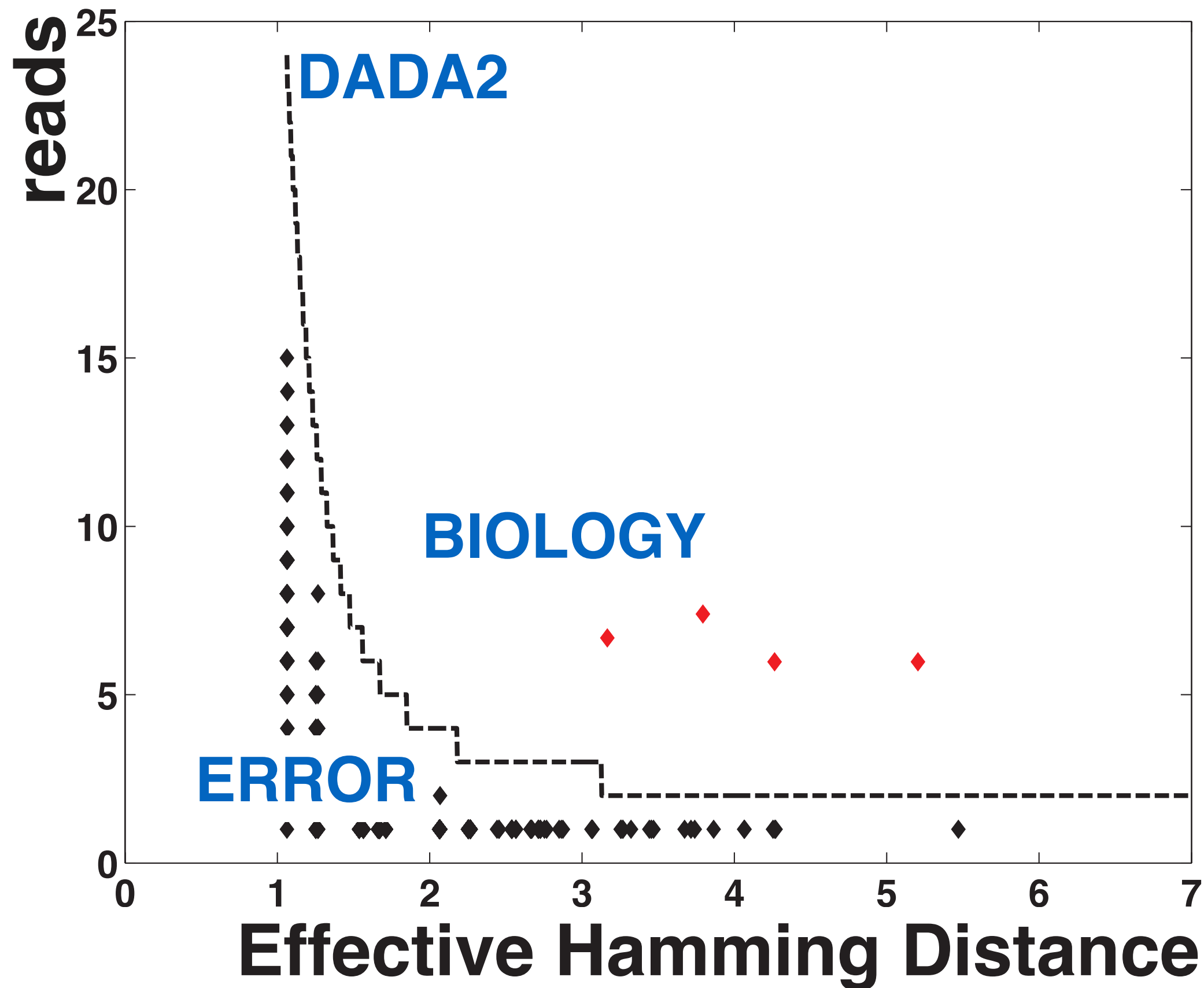


# Signal from Noise: OTUs





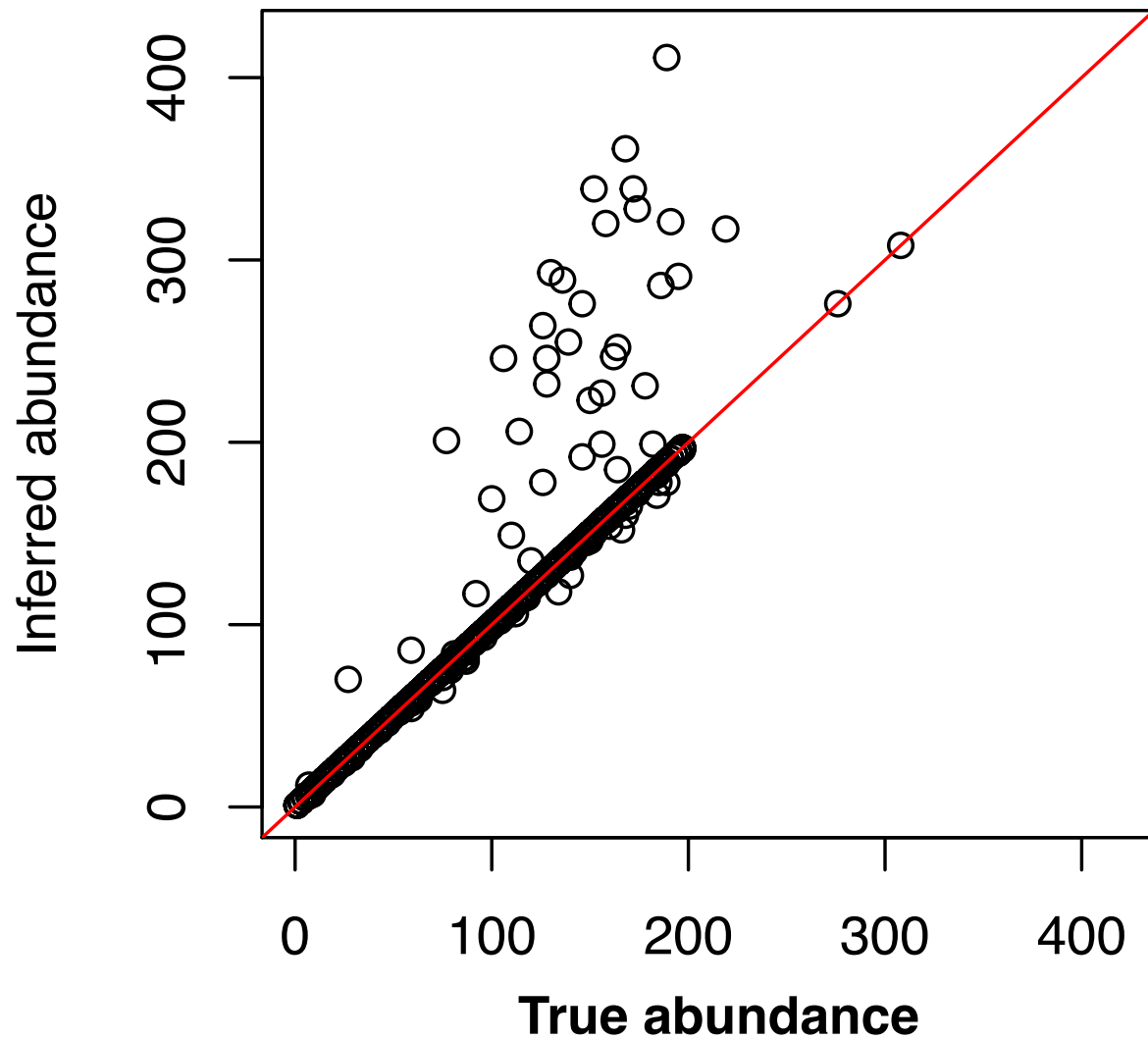
# Signal from Noise: DADA2



# **Accuracy** and **Resolution**

# Accuracy: Simulated data

mothur (an)



**TP:** 978

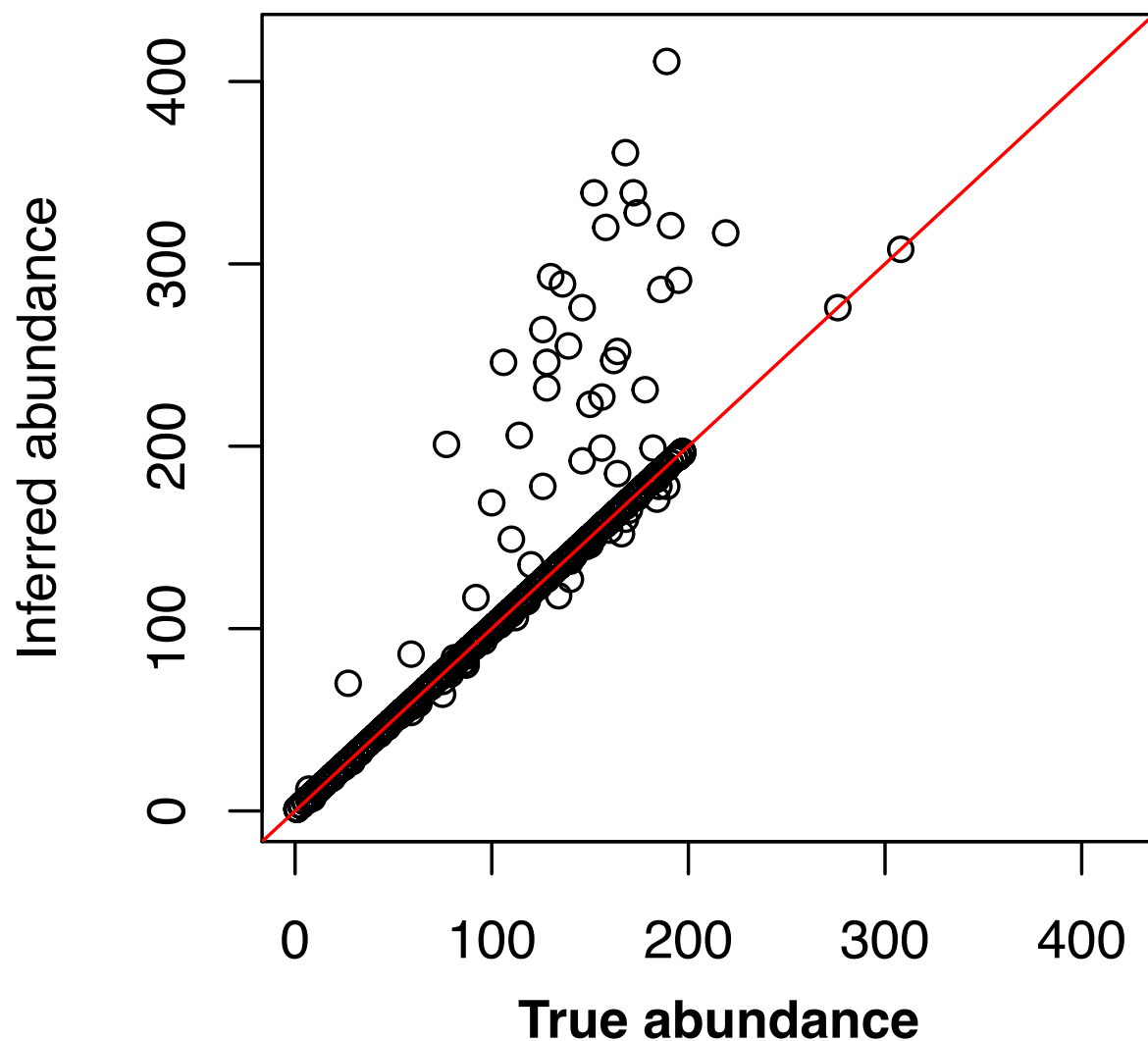
**FP:** 272

**FN:** 77

**cor:** 0.935

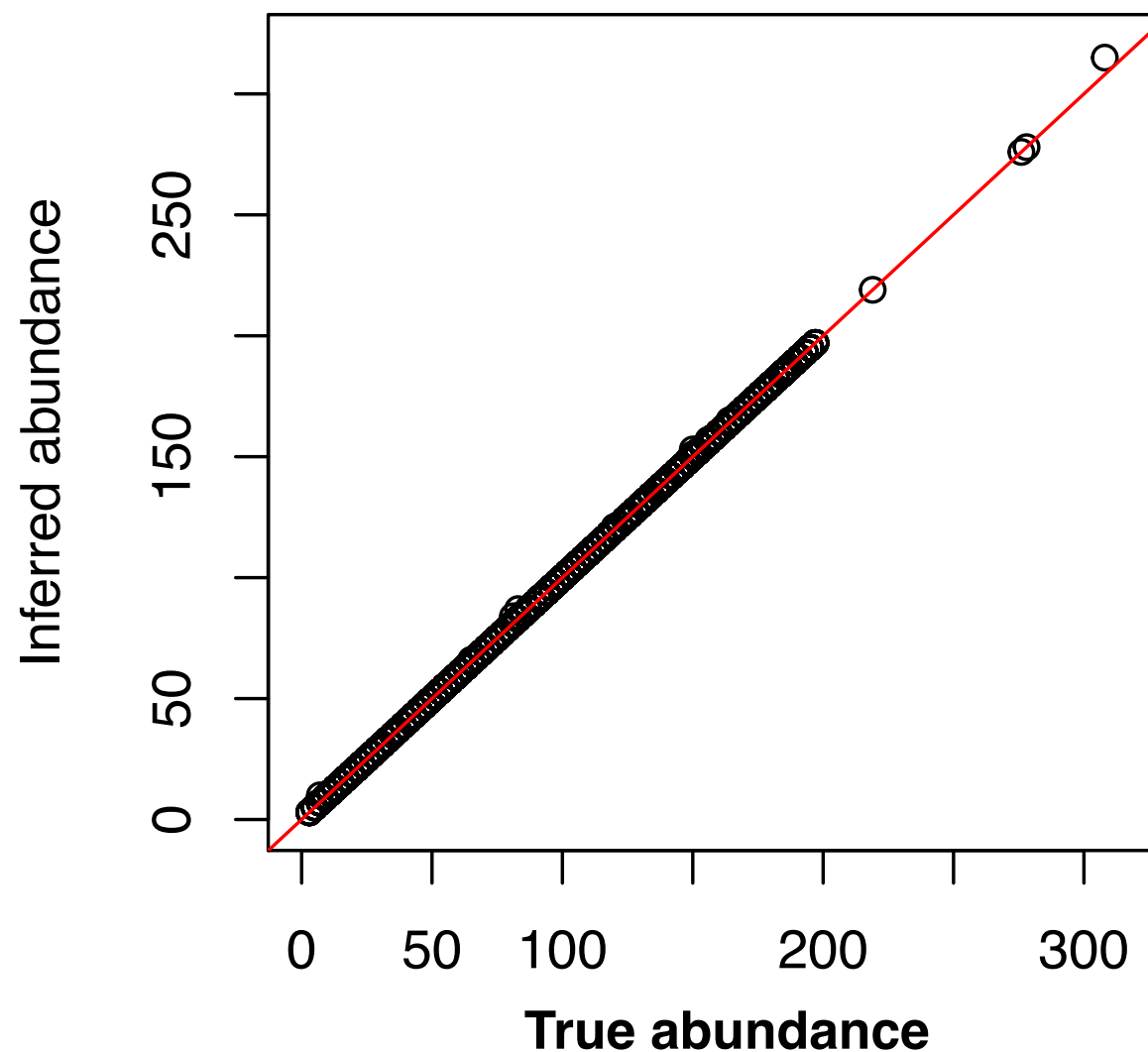
# Accuracy: Simulated data

mothur (an)



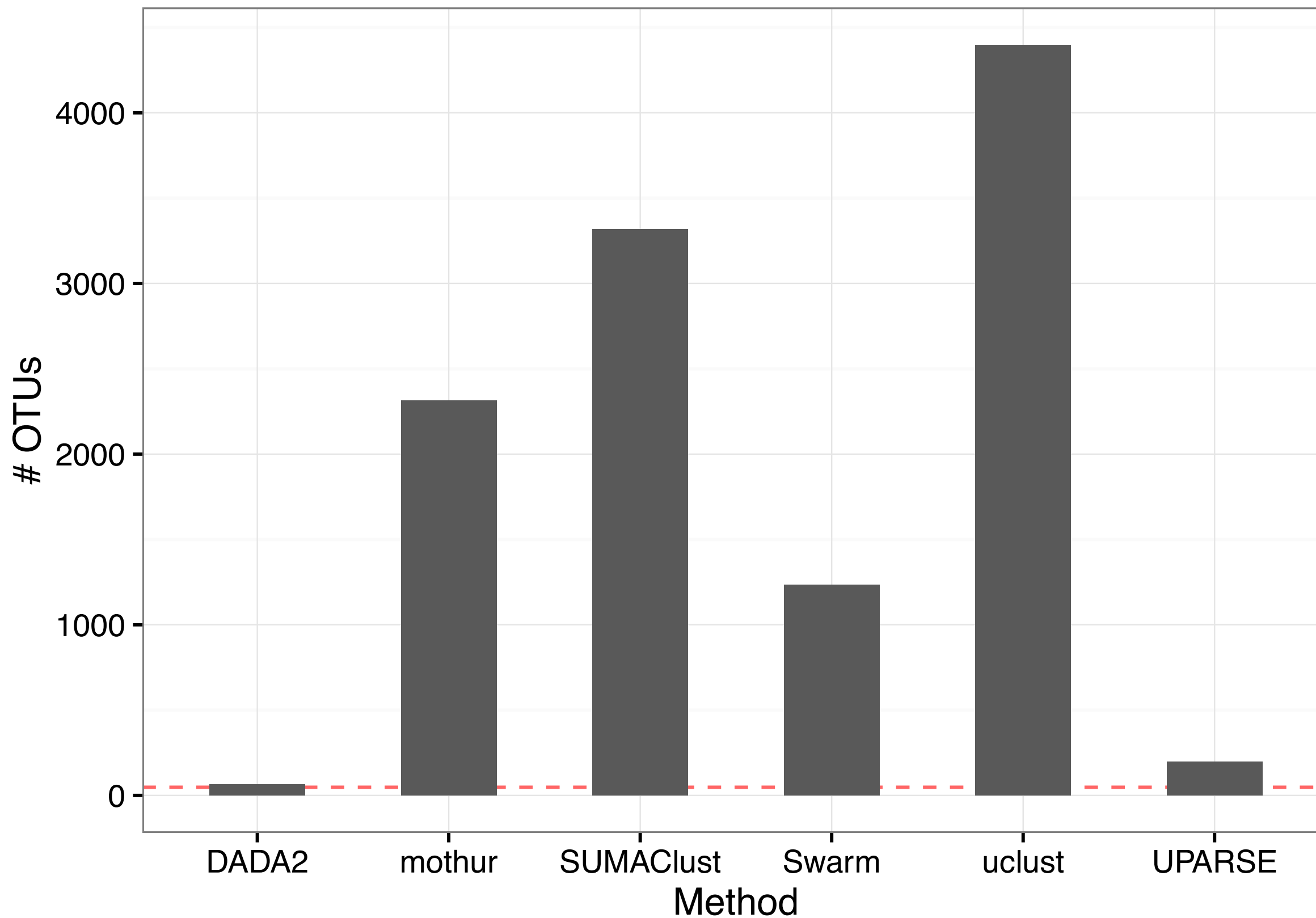
**TP:** 978  
**FP:** 272  
**FN:** 77  
**cor:** 0.935

DADA2



**TP:** 1042  
**FP:** 0  
**FN:** 13  
**cor:** 0.999

# Accuracy: Mock community



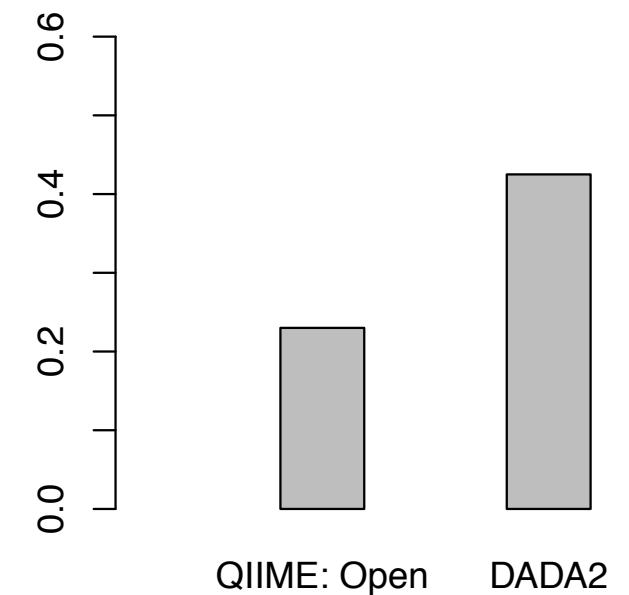
**Credit:** Kopylova, et al. mSystems, 2016.

# Accuracy: Arsenic treatment

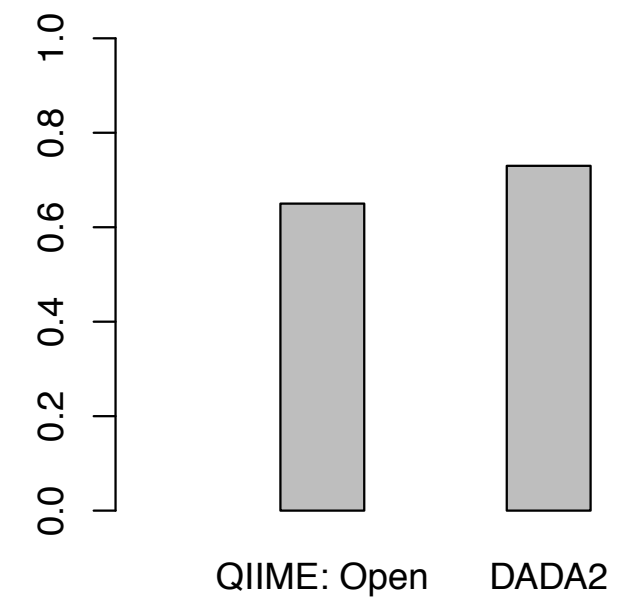
**Unpublished  
ordination  
removed.**

## Variance explained by

PC1 + PC2

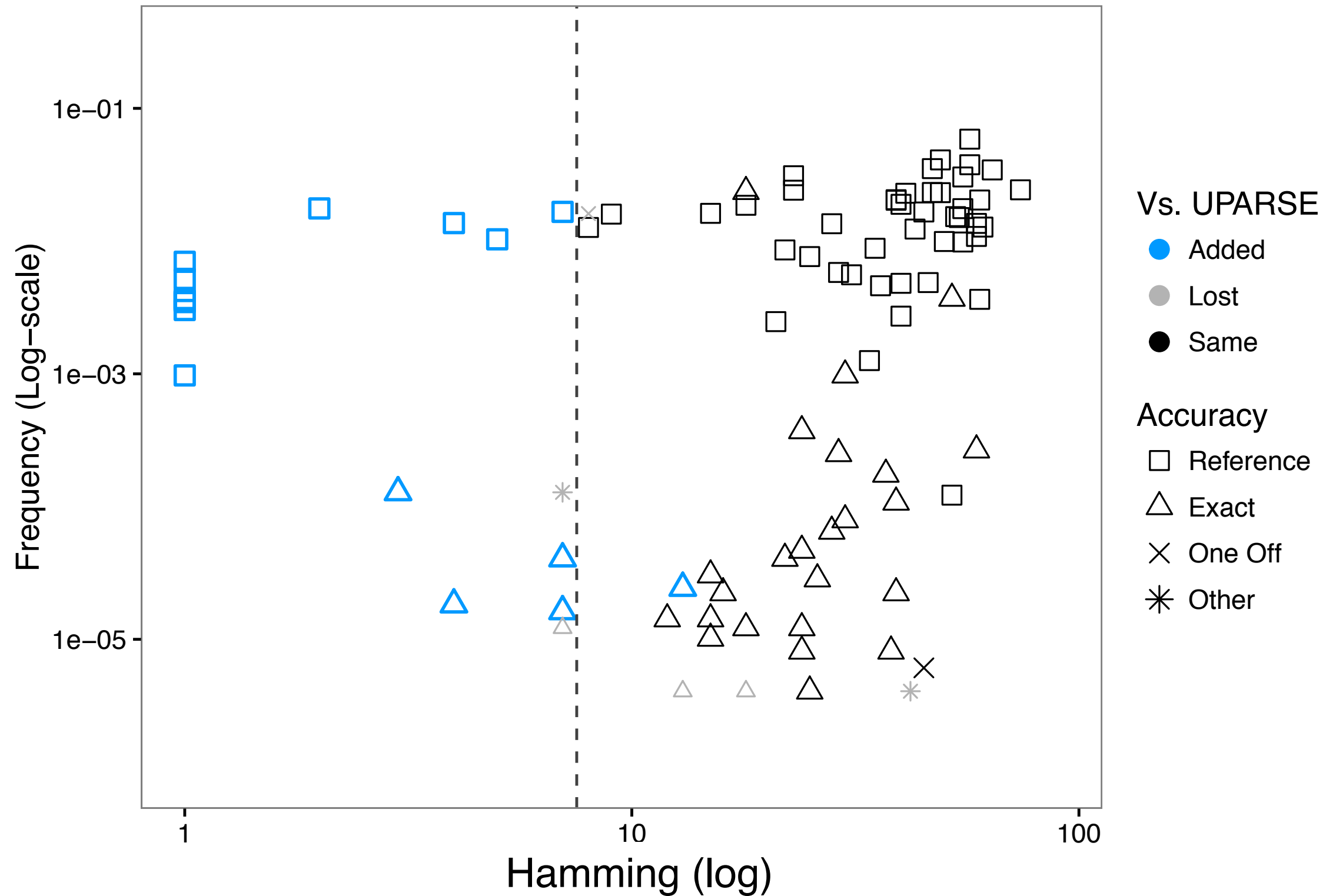


**Arsenic**



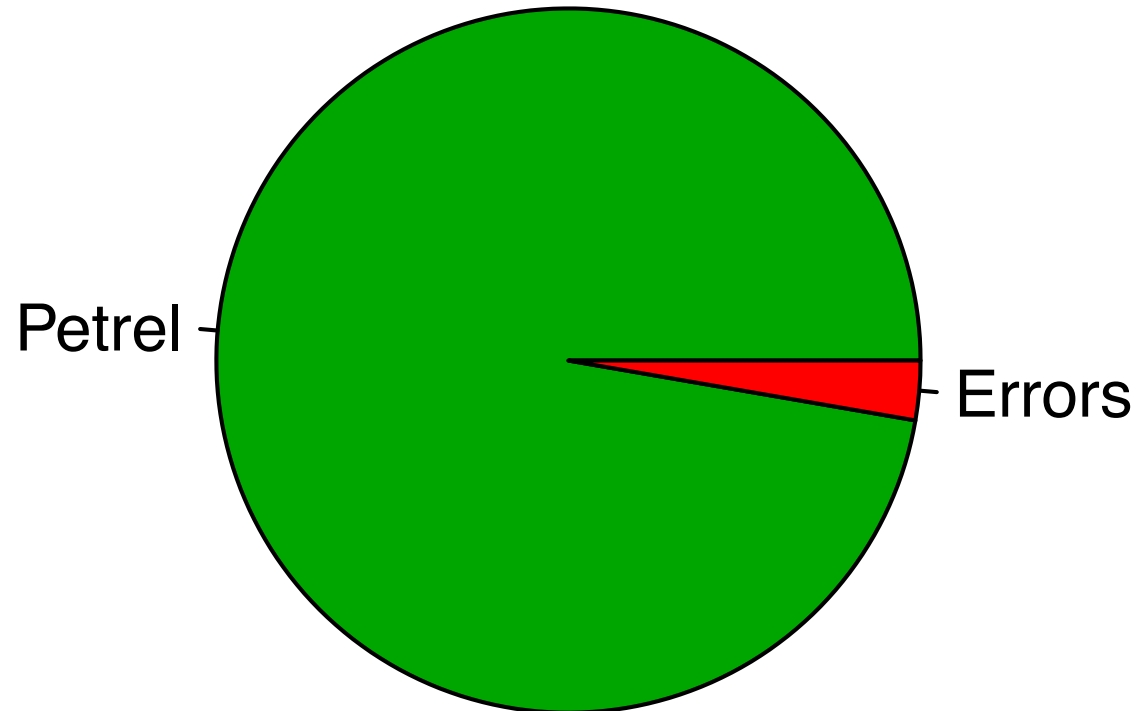
# Resolution: Mock Community

## DADA2

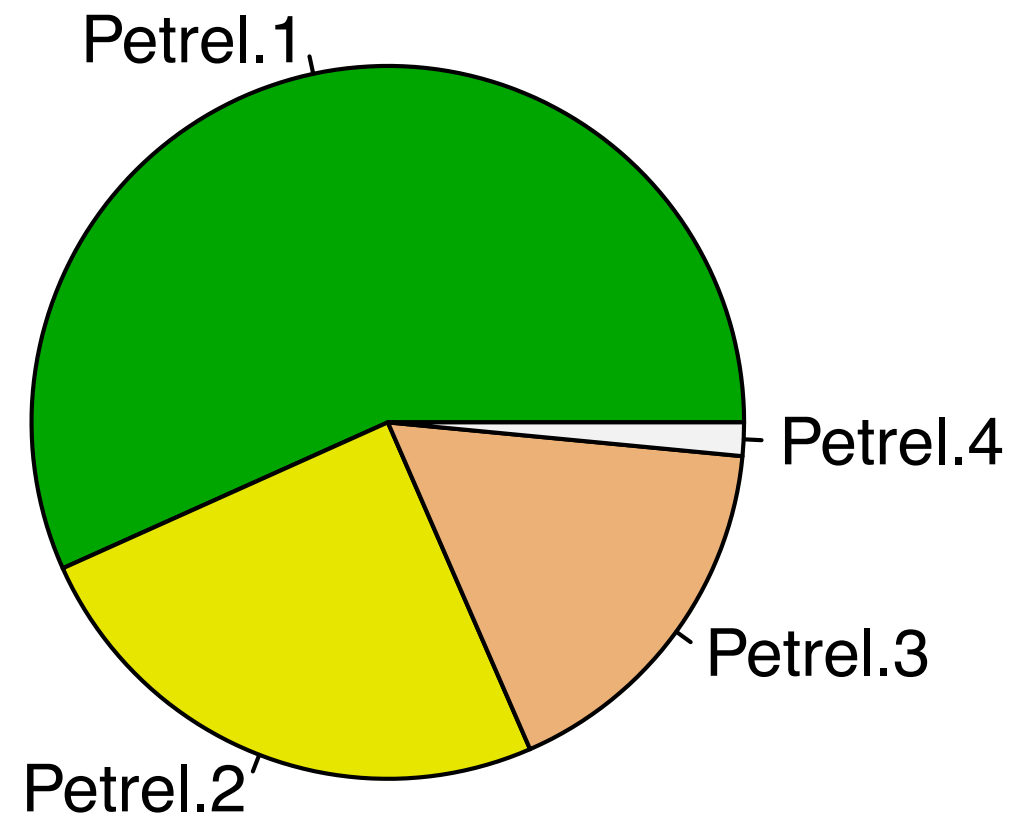


# Resolution: Petrel aDNA

**QIIME: De novo**

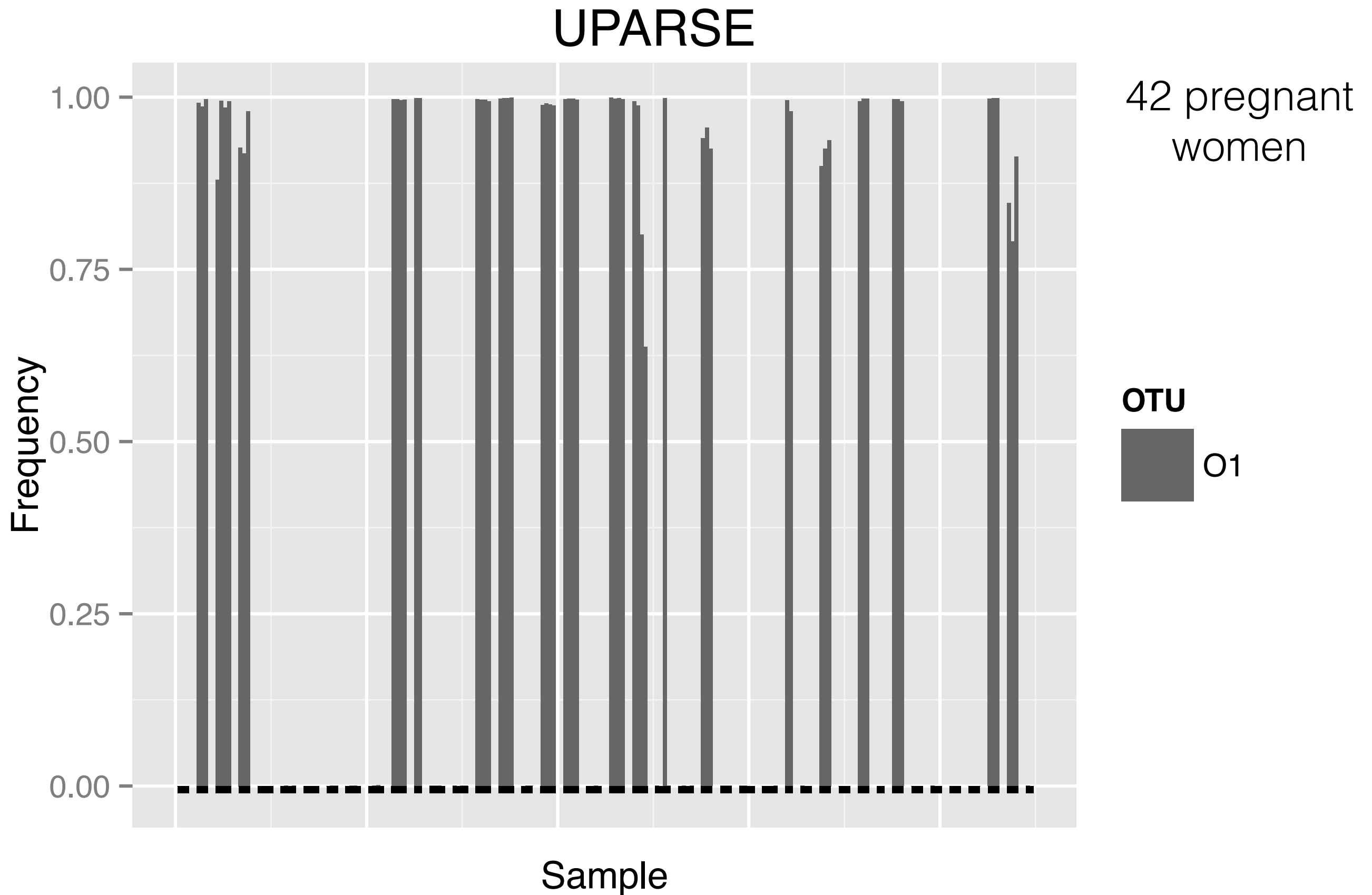


**DADA2**





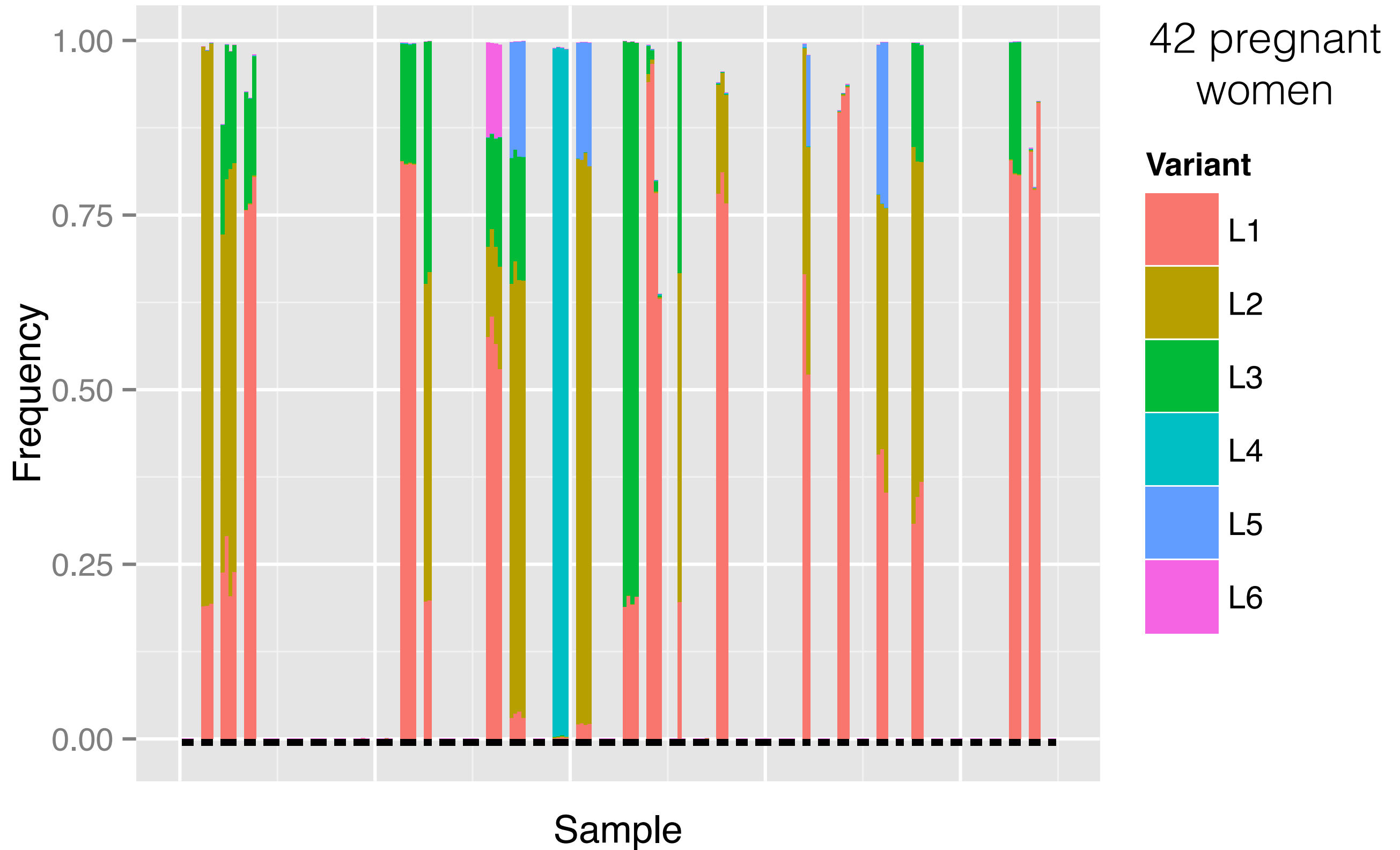
# Resolution: *L. crispatus*



**Data:** MacIntyre et al. Scientific Reports, 2015.

# Resolution: *L. crispatus*

## DADA2



**Data:** MacIntyre et al. Scientific Reports, 2015.

# **Scaling and Comparison**

# The sequence is the label

ATTAAACGAGATTATAACCCAGAGTACGAATA...

is a *consistent label*

OTU85 is *not*

# The sequence is the label

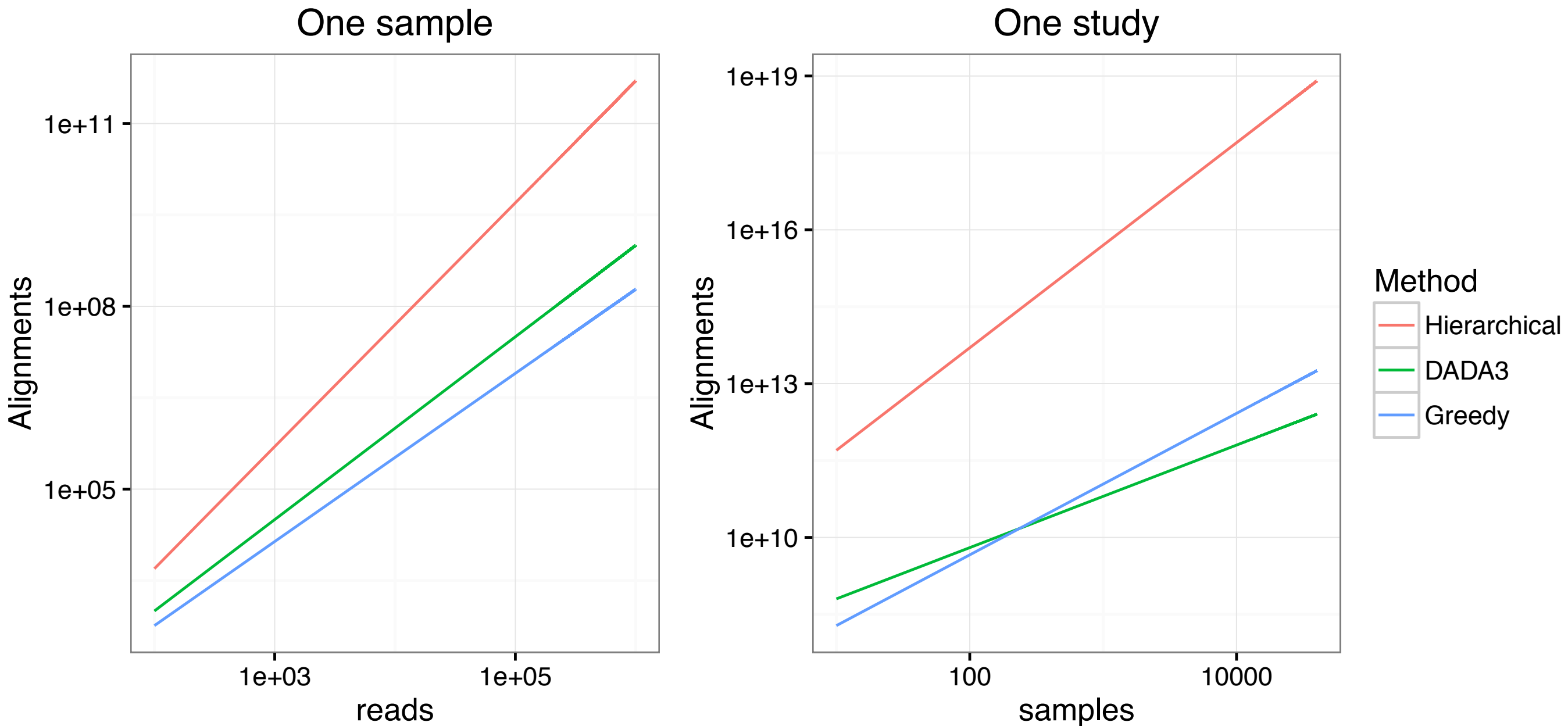
ATTAAACGAGATTATAACCCAGAGTACGAATA...

is a *consistent label*

OTU85 is *not*

Consistent labels  
allow samples to be  
**independently processed**

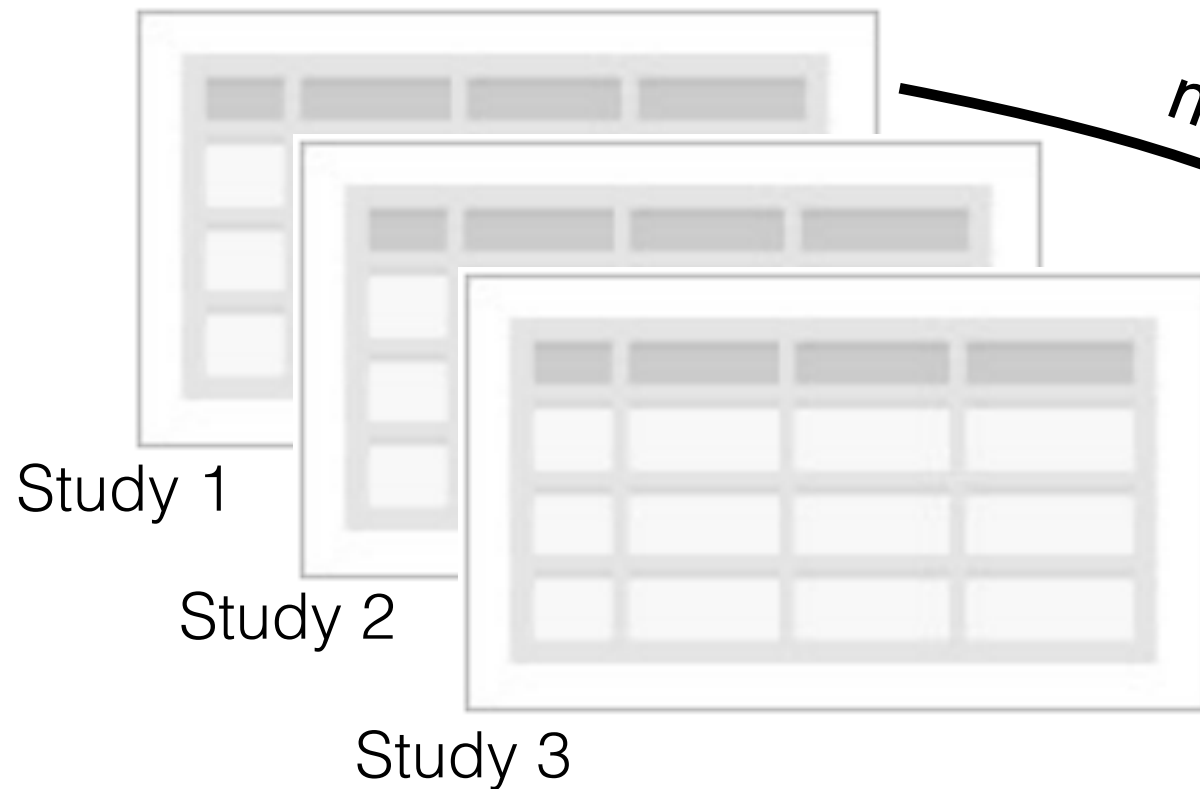
# Consistent Labels: Scaling



Separable processing.  
Flat memory requirements.  
Task parallelization.

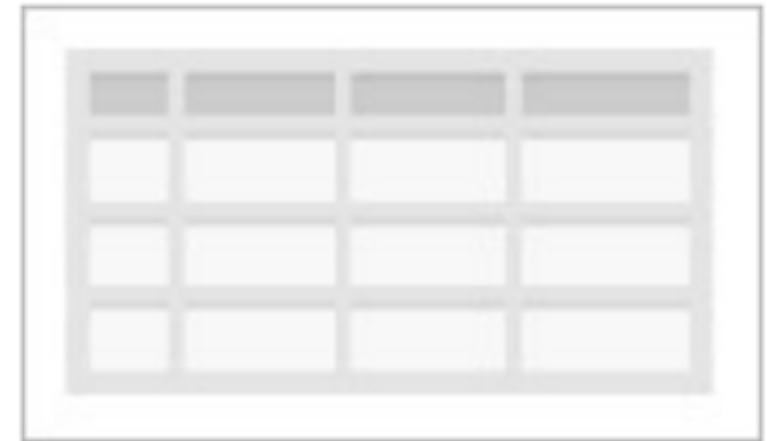
# Consistent Labels: Comparison

Sequence Tables



merge

Cross-study comparison



Eliminates need for joint reprocessing of raw data.

# DADA2 is...

- A replacement for OTU picking
- Accurate and high-resolution
- Implemented in an R package
- Open source





# DADA2 is...

- A replacement for OTU picking
- Accurate and high-resolution
- Implemented in an R package
- Open source



*NATURE METHODS* | BRIEF COMMUNICATION

## DADA2: High-resolution sample inference from Illumina amplicon data

[Benjamin J Callahan](#), [Paul J McMurdie](#), [Michael J Rosen](#), [Andrew W Han](#), [Amy Jo A Johnson](#) & [Susan P Holmes](#)

# Acknowledgements



Susan Holmes



Joey McMurdie



Michael Rosen



National Institutes of Health