

The DADA2 Method

The amplicon inference problem

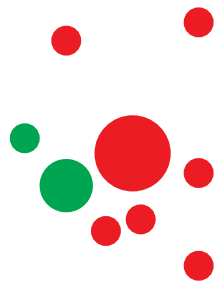
Infer the sample types and abundances $\{(s, a)\}$
from error-ful amplicon reads $\{r\}$.

sample
sequences

amplicon reads



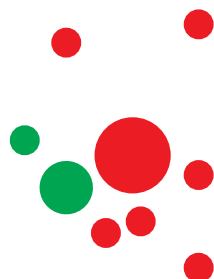
Errors



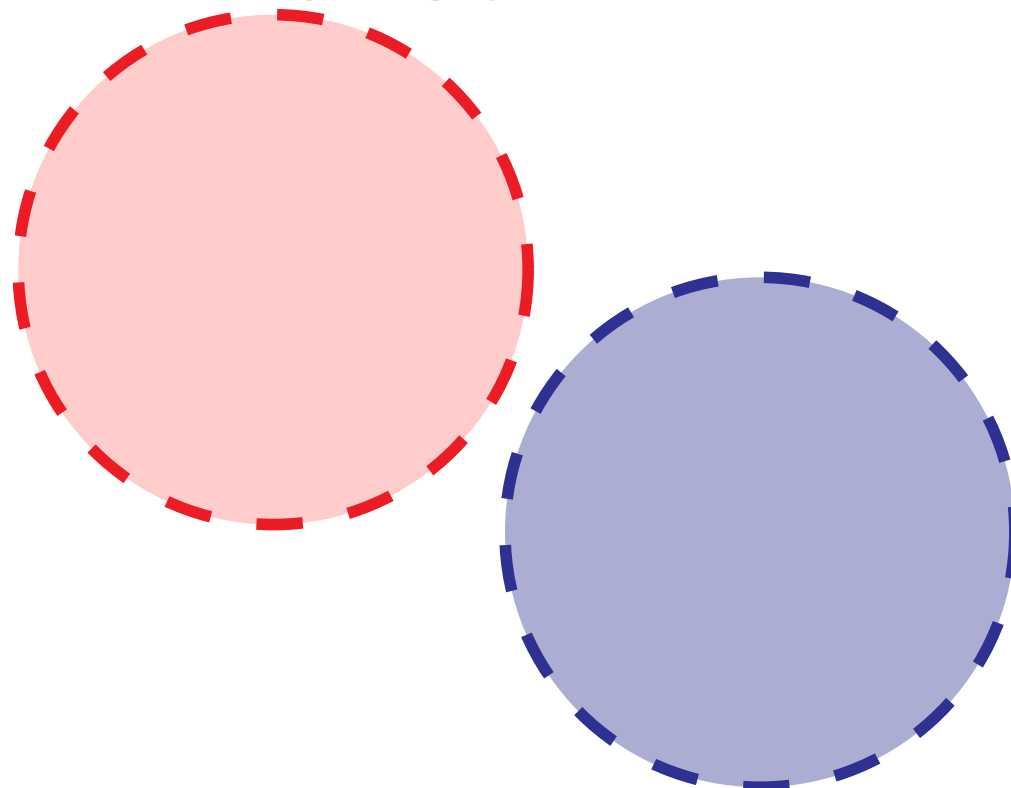
sample
sequences



amplicon reads



OTUs

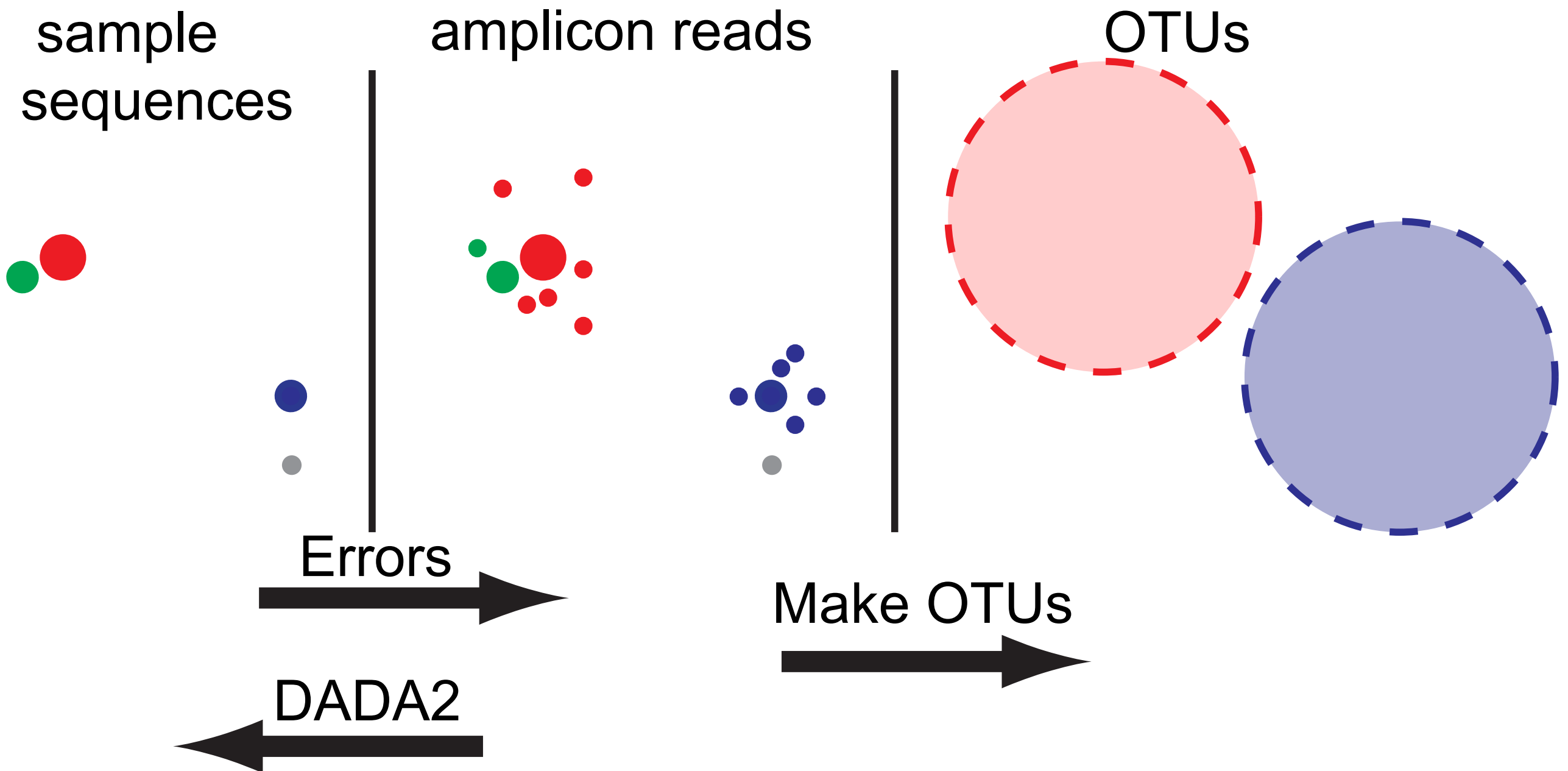


Errors



Make OTUs







Amplicon Sequencing. **Exactly.** ***Version 1.10***

Error Model

An Error Model

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

An Error Model

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

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

An Error Model

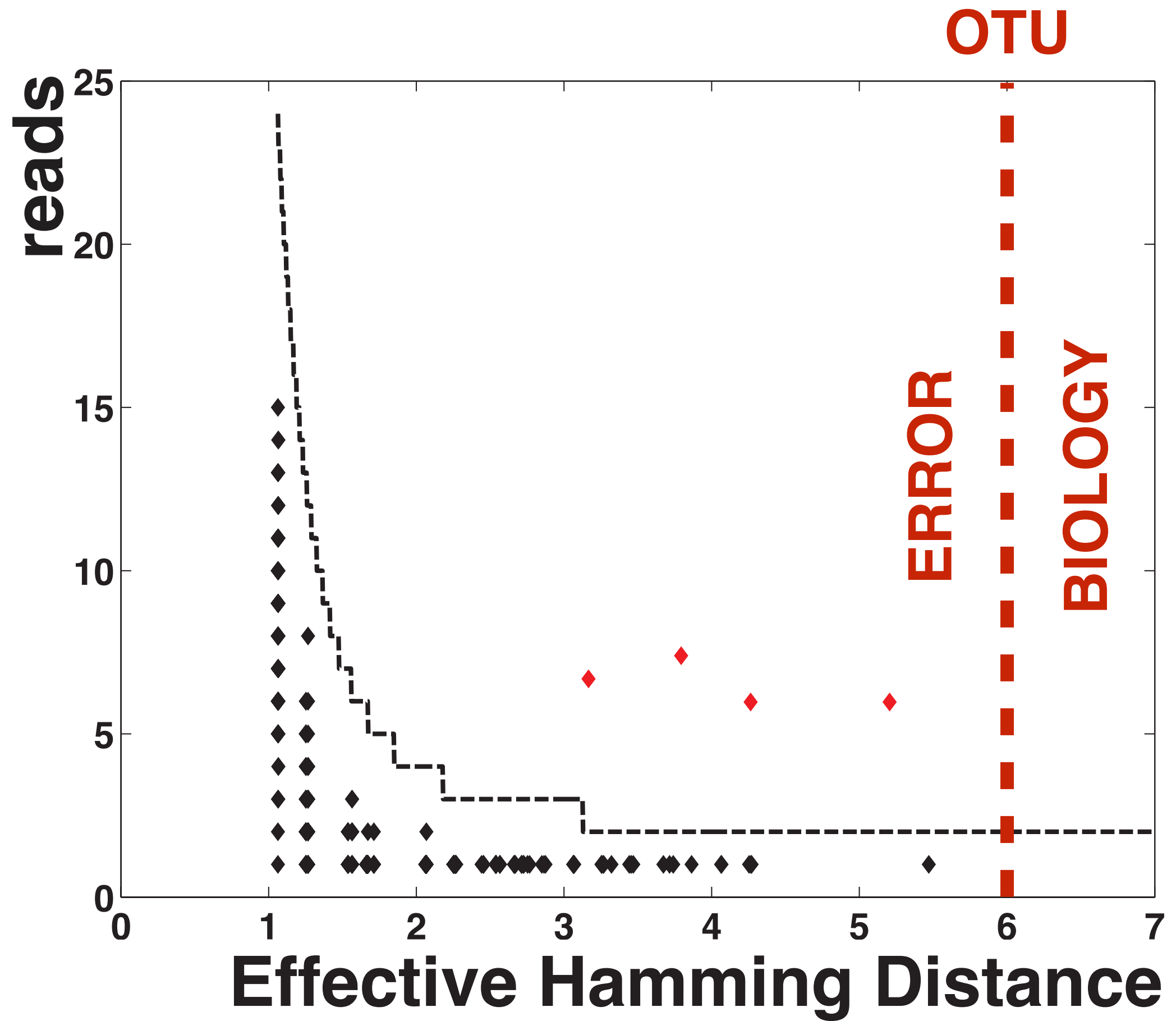
s : ATTAACGAGATTATAACCCAGAGTACGAATA . . .
 | |
r : ATCAACGAGATTATAACAAGAGTACGAATA . . .

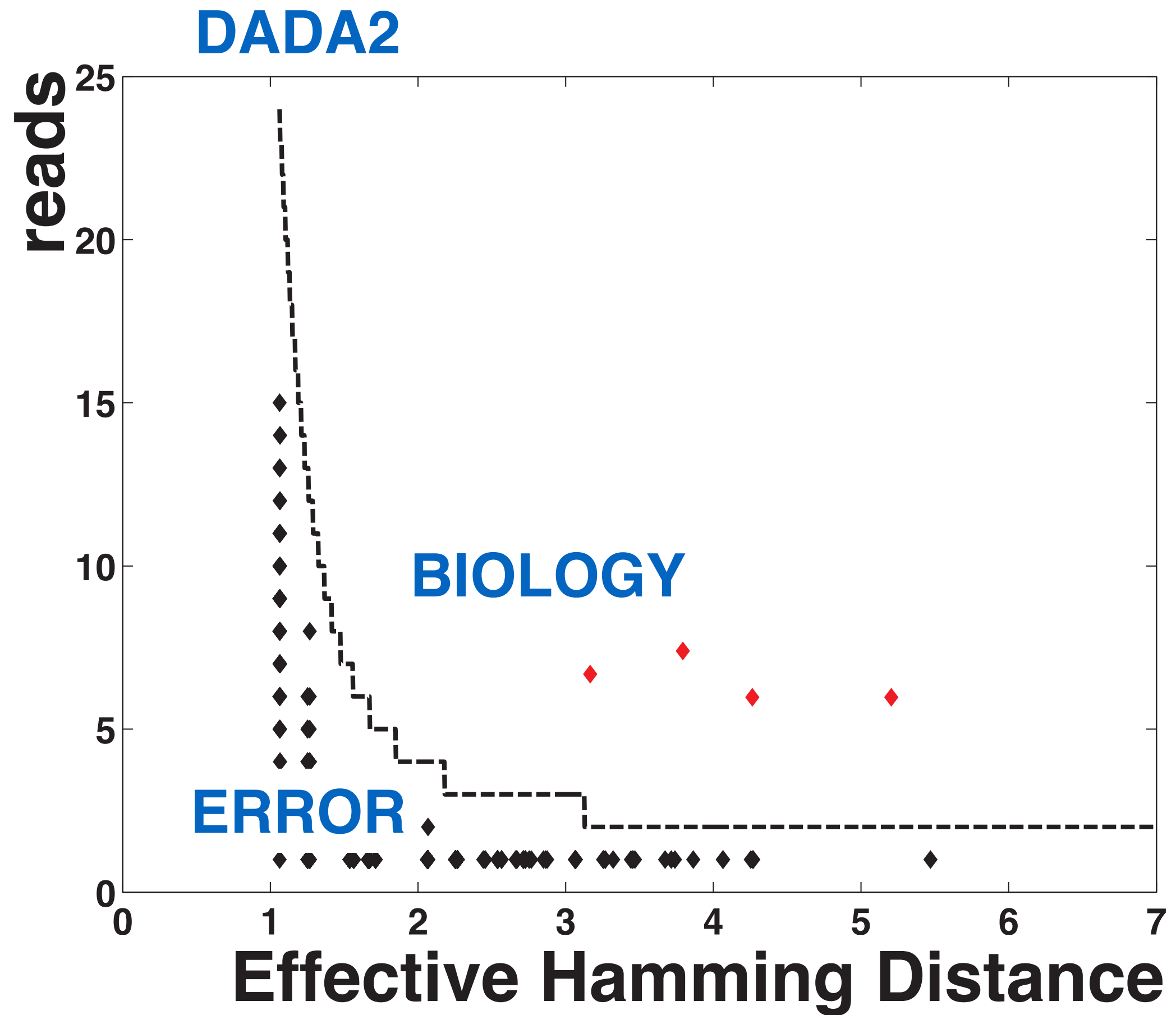
$$p(r|s) = \prod_{i=1}^L p(r(i)|s(i), q_r(i), Z)$$

Error process is independent across nucleotides.

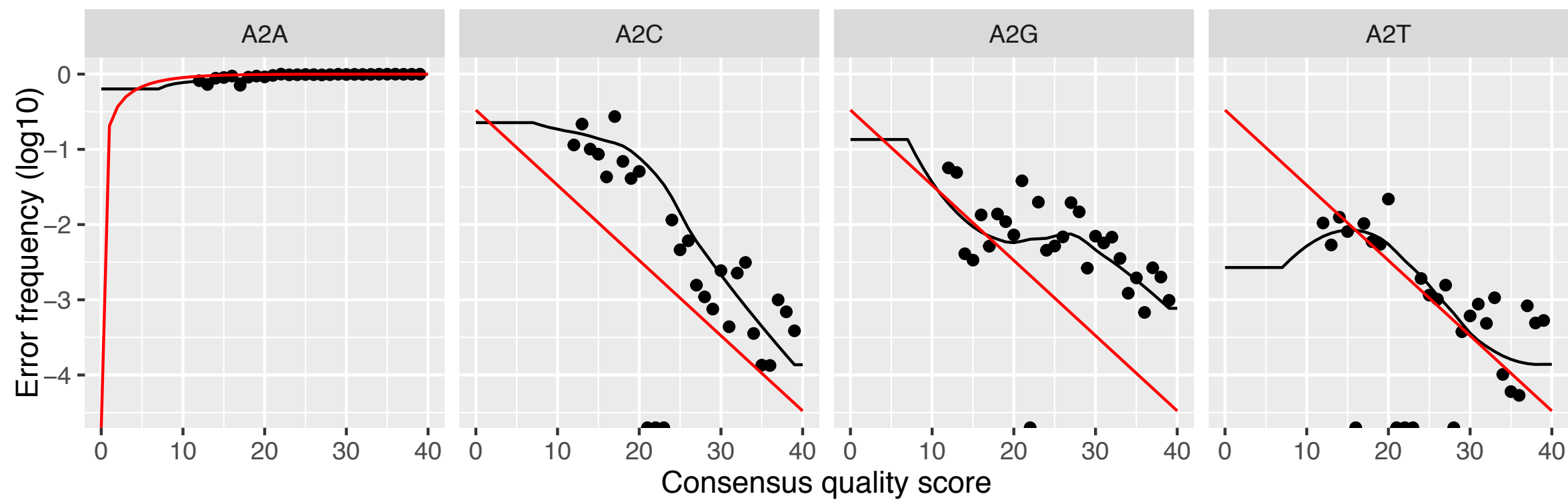
Per-nucleotide transition rate depends on:

- Sample nucleotide
- Read nucleotide
- Read quality at that position
- Batch effect (eg. run)

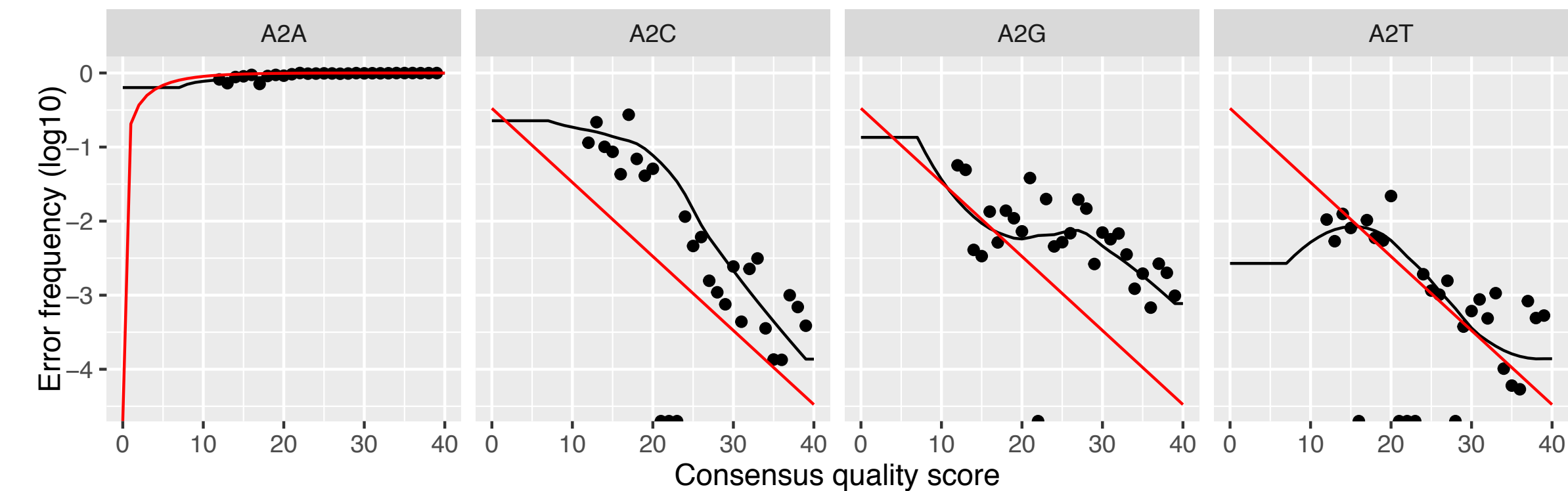




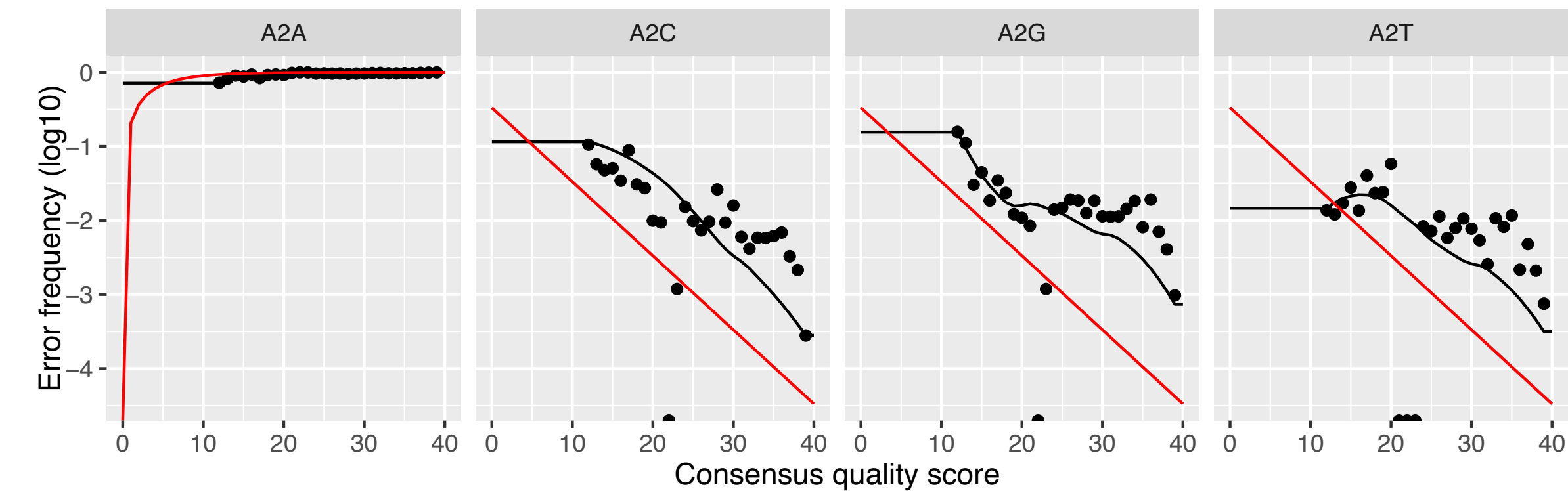
Learning Errors



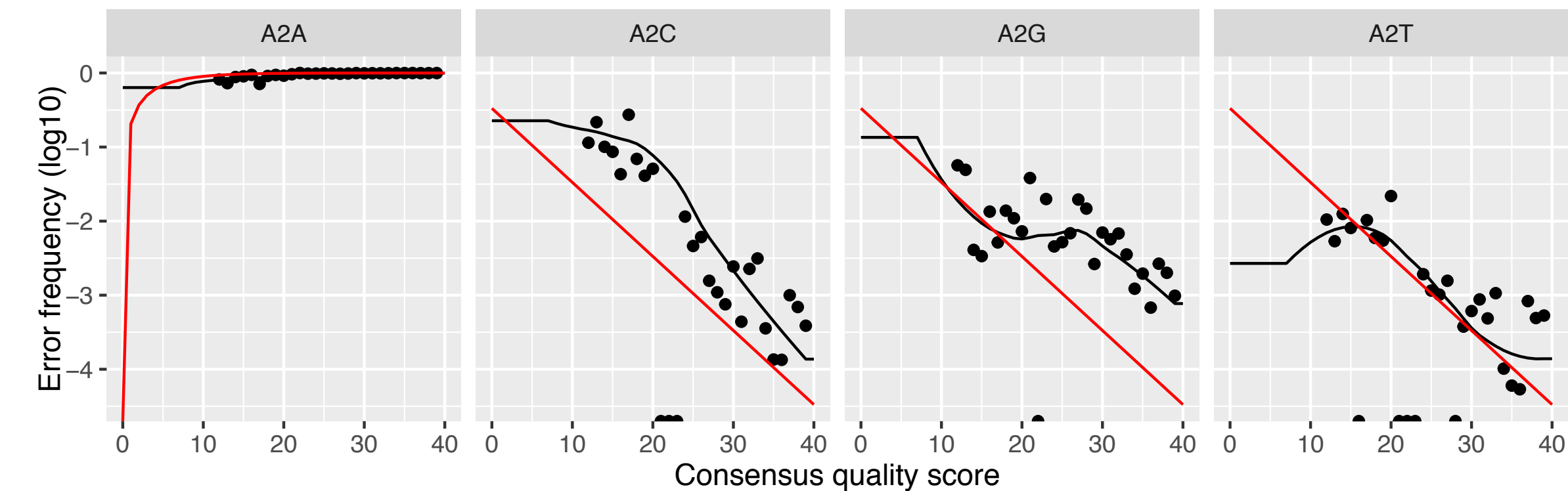
Study A



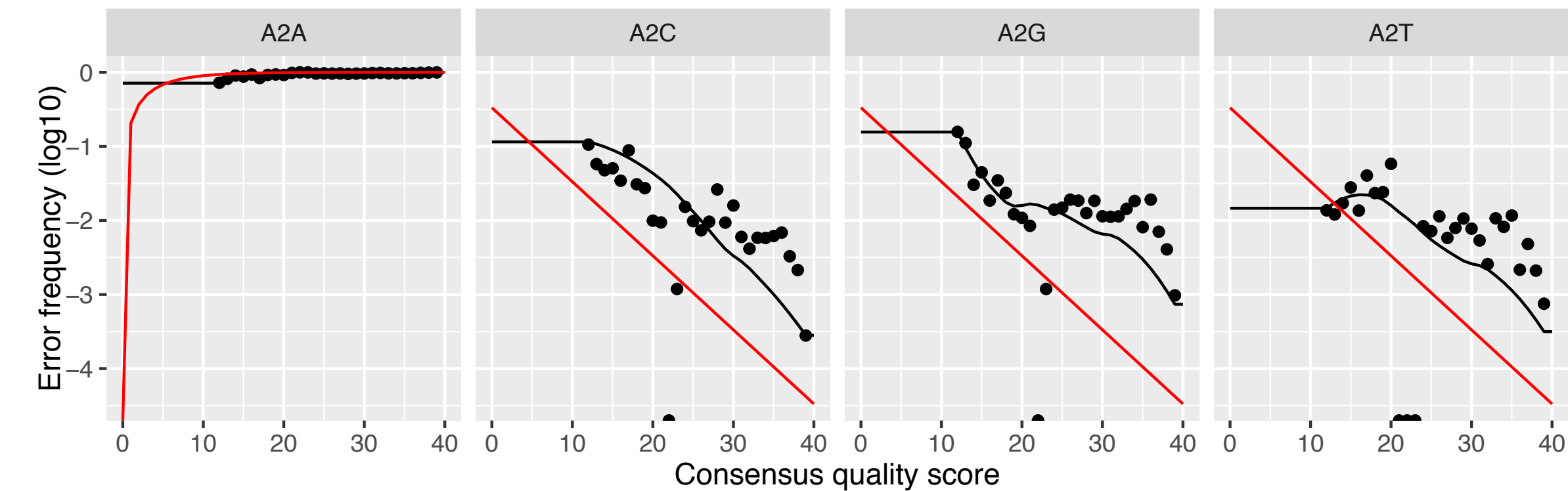
Study B



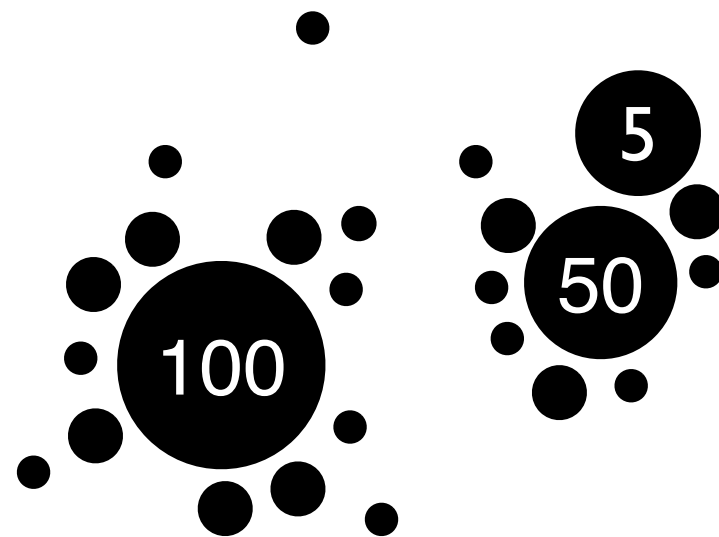
Study A



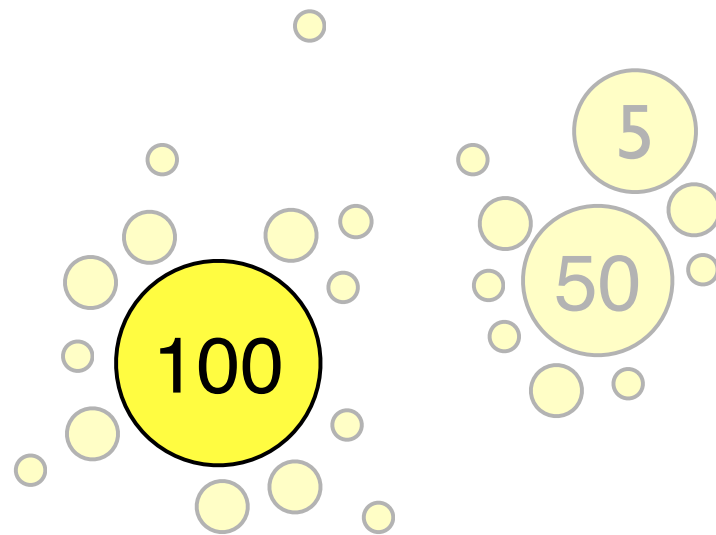
Study B



But How?



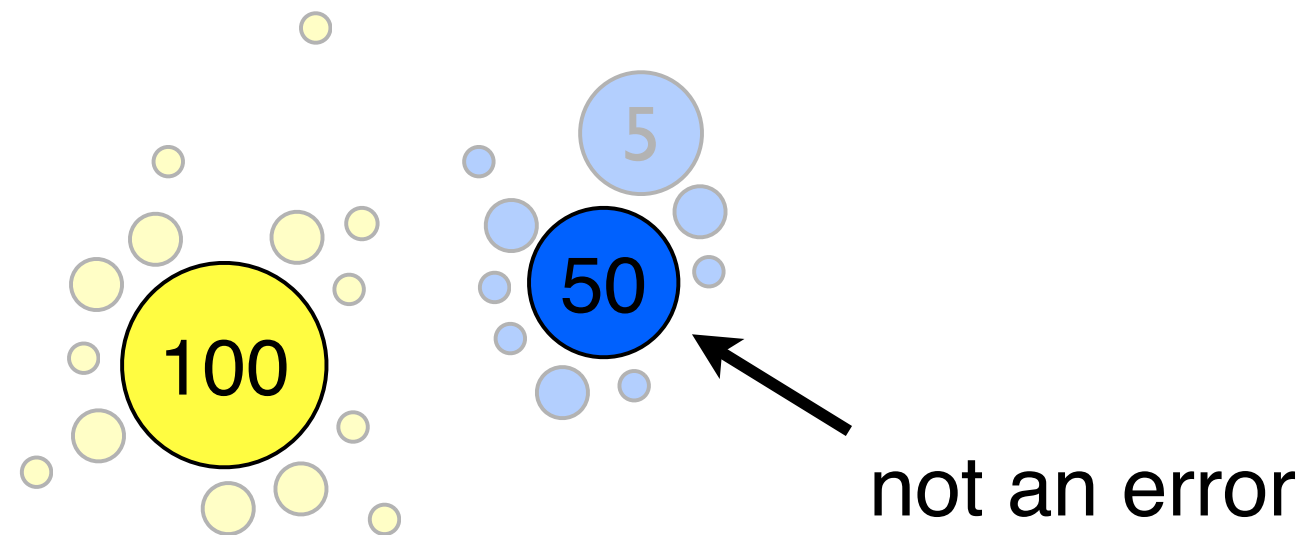
Initial guess: one real sequence + errors



Infer initial *error model* under this assumption.

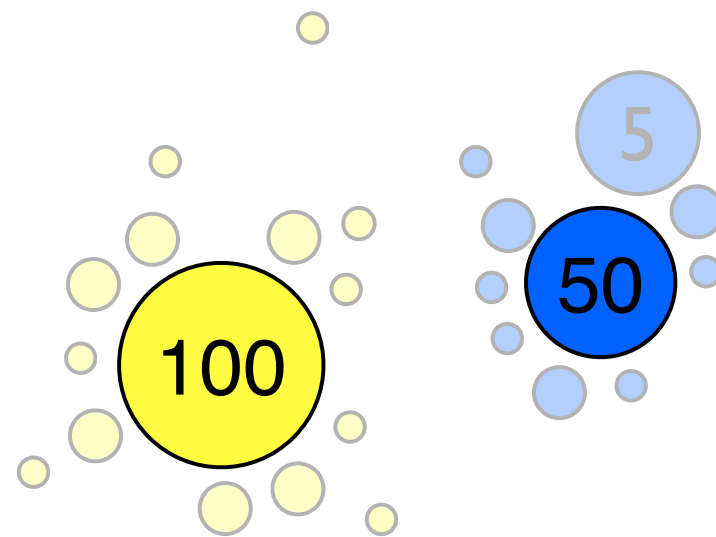
$\text{Pr}(i \rightarrow j) =$

| | A | C | G | T |
|---|-----------|-----------|-----------|-----------|
| A | 0.97 | 10^{-2} | 10^{-2} | 10^{-2} |
| C | 10^{-2} | 0.97 | 10^{-2} | 10^{-2} |
| G | 10^{-2} | 10^{-2} | 0.97 | 10^{-2} |
| T | 10^{-2} | 10^{-2} | 10^{-2} | 0.97 |



Reject unlikely error under model. **Recruit** errors.

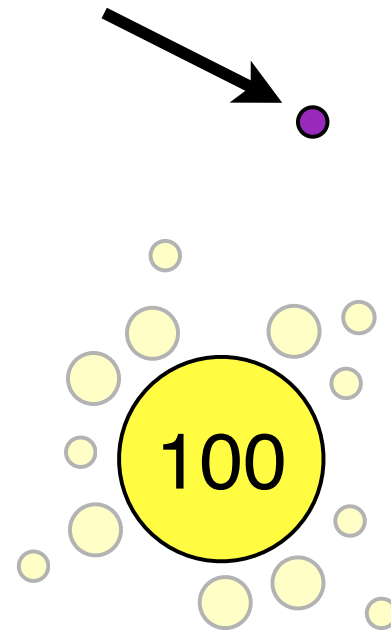
| | A | C | G | T |
|---|-----------|-----------|-----------|-----------|
| A | 0.97 | 10^{-2} | 10^{-2} | 10^{-2} |
| C | 10^{-2} | 0.97 | 10^{-2} | 10^{-2} |
| G | 10^{-2} | 10^{-2} | 0.97 | 10^{-2} |
| T | 10^{-2} | 10^{-2} | 10^{-2} | 0.97 |



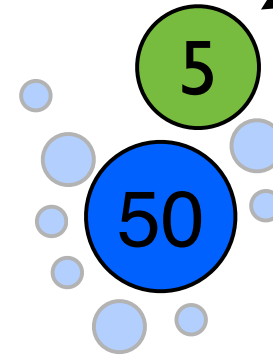
Update the model.

| | A | C | G | T |
|---|-----------|-----------|-----------|-----------|
| A | 0.997 | 10^{-3} | 10^{-3} | 10^{-3} |
| C | 10^{-3} | 0.997 | 10^{-3} | 10^{-3} |
| G | 10^{-3} | 10^{-3} | 0.997 | 10^{-3} |
| T | 10^{-3} | 10^{-3} | 10^{-3} | 0.997 |

not an error

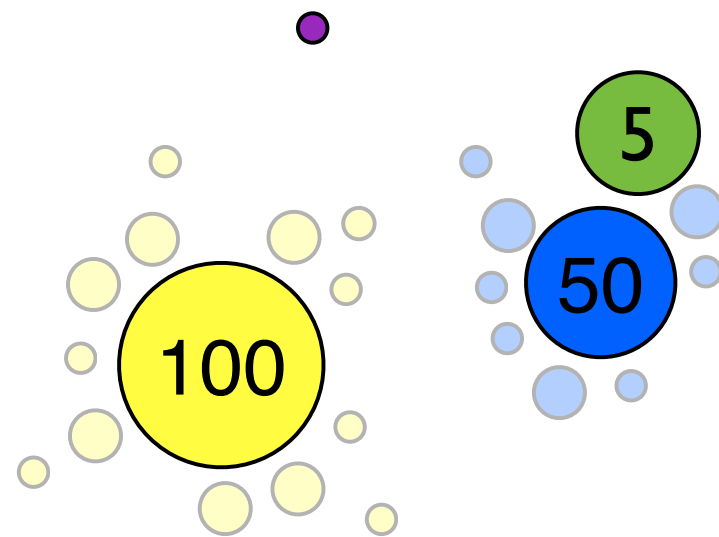


not an error



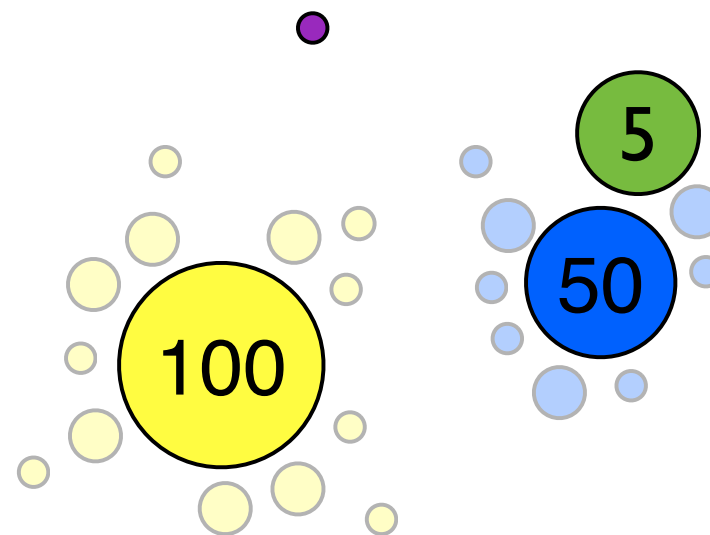
Reject more sequences under *new* model

| | A | C | G | T |
|---|-----------|-----------|-----------|-----------|
| A | 0.997 | 10^{-3} | 10^{-3} | 10^{-3} |
| C | 10^{-3} | 0.997 | 10^{-3} | 10^{-3} |
| G | 10^{-3} | 10^{-3} | 0.997 | 10^{-3} |
| T | 10^{-3} | 10^{-3} | 10^{-3} | 0.997 |



Update model again

| | A | C | G | T |
|---|--------------------|--------------------|--------------------|--------------------|
| A | 0.998 | 1×10^{-4} | 2×10^{-3} | 2×10^{-4} |
| C | 6×10^{-5} | 0.999 | 3×10^{-6} | 1×10^{-3} |
| G | 1×10^{-3} | 3×10^{-6} | 0.999 | 6×10^{-5} |
| T | 2×10^{-4} | 2×10^{-3} | 1×10^{-4} | 0.998 |



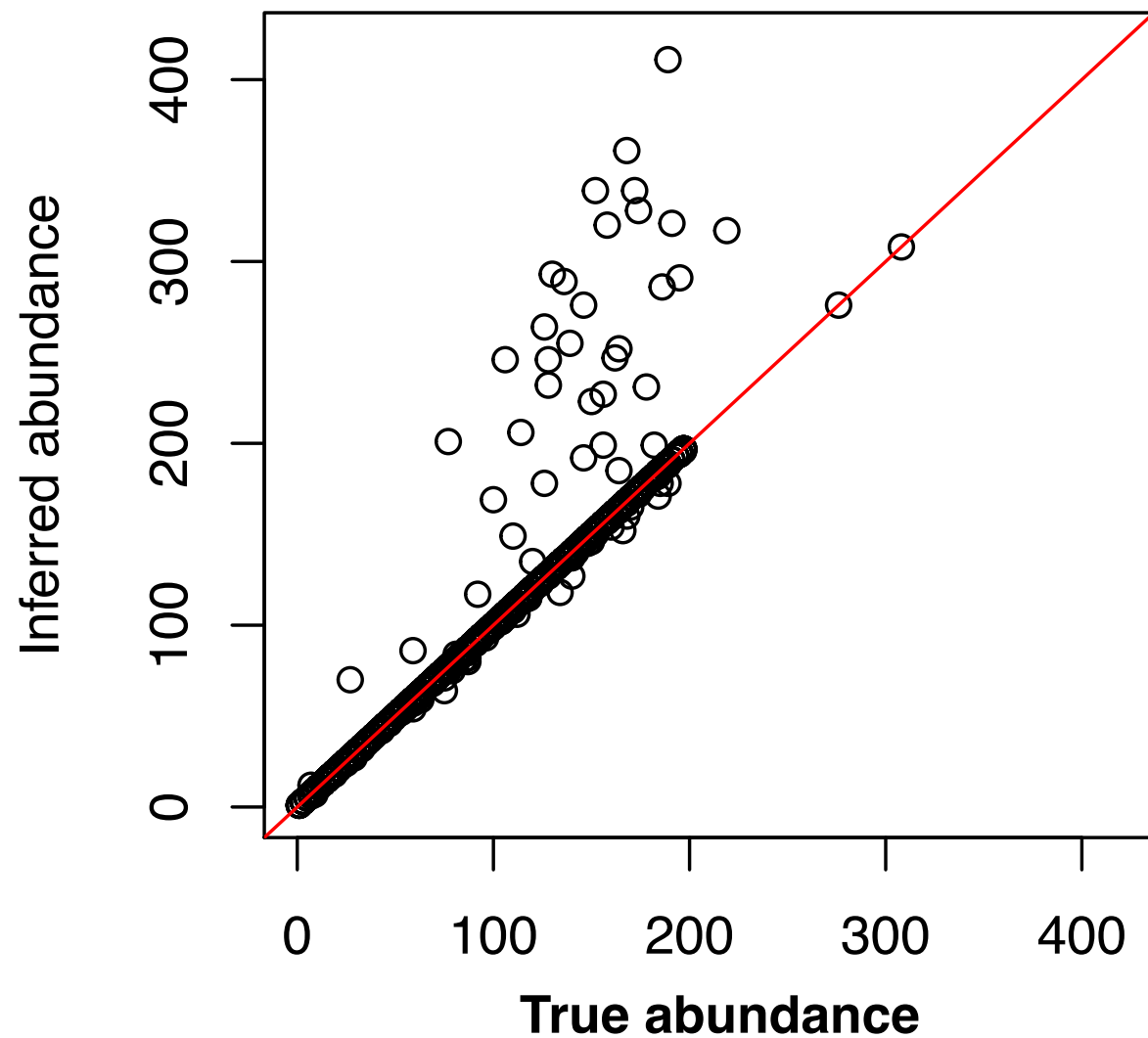
Convergence: all errors are plausible

| | A | C | G | T |
|---|--------------------|--------------------|--------------------|--------------------|
| A | 0.998 | 1×10^{-4} | 2×10^{-3} | 2×10^{-4} |
| C | 6×10^{-5} | 0.999 | 3×10^{-6} | 1×10^{-3} |
| G | 1×10^{-3} | 3×10^{-6} | 0.999 | 6×10^{-5} |
| T | 2×10^{-4} | 2×10^{-3} | 1×10^{-4} | 0.998 |

Accuracy and **Resolution**

Accuracy: Simulated data

mothur (an)



TP: 978

FP: 272

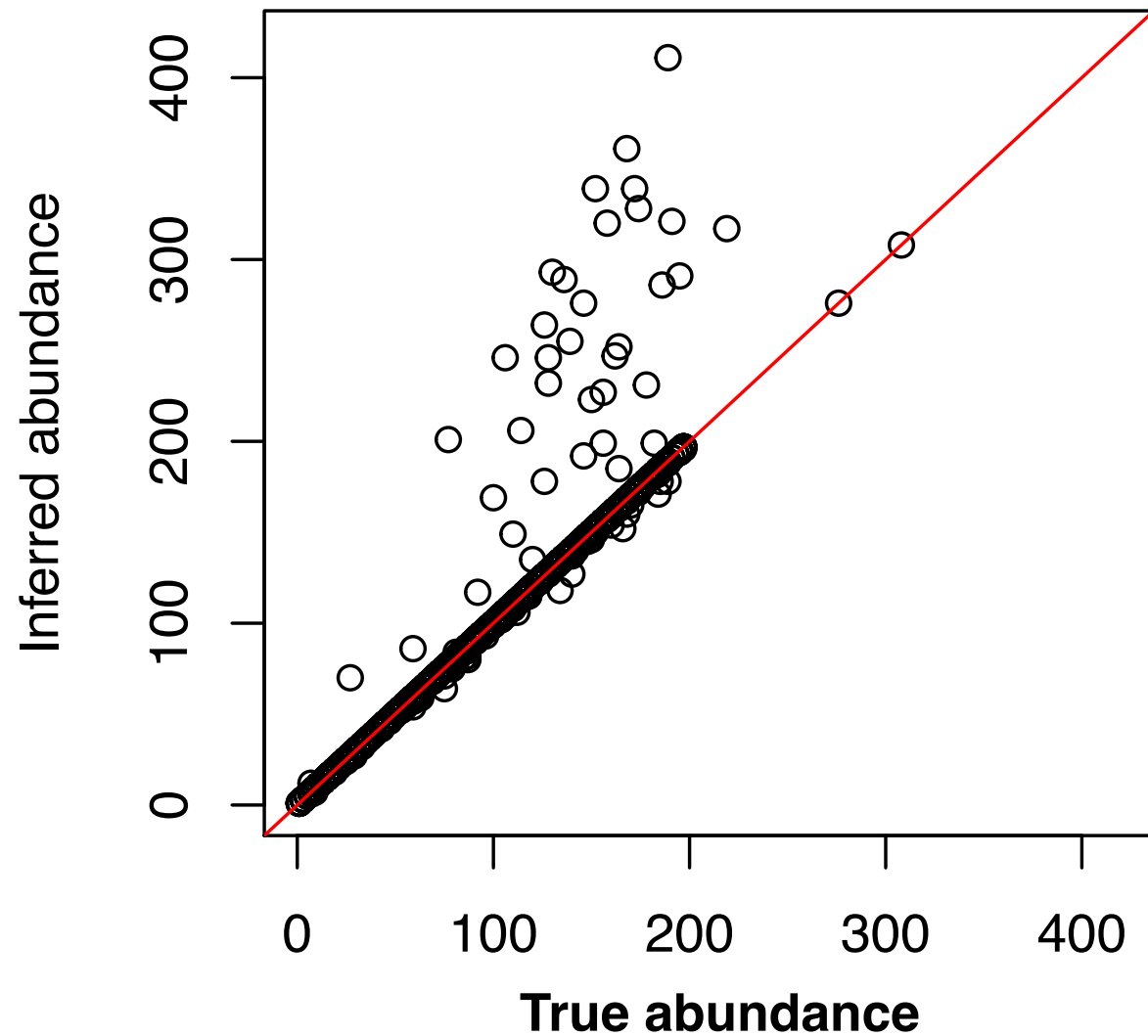
FN: 77

cor: 0.935

Data: Kopylova, et al. mSystems, 2016.

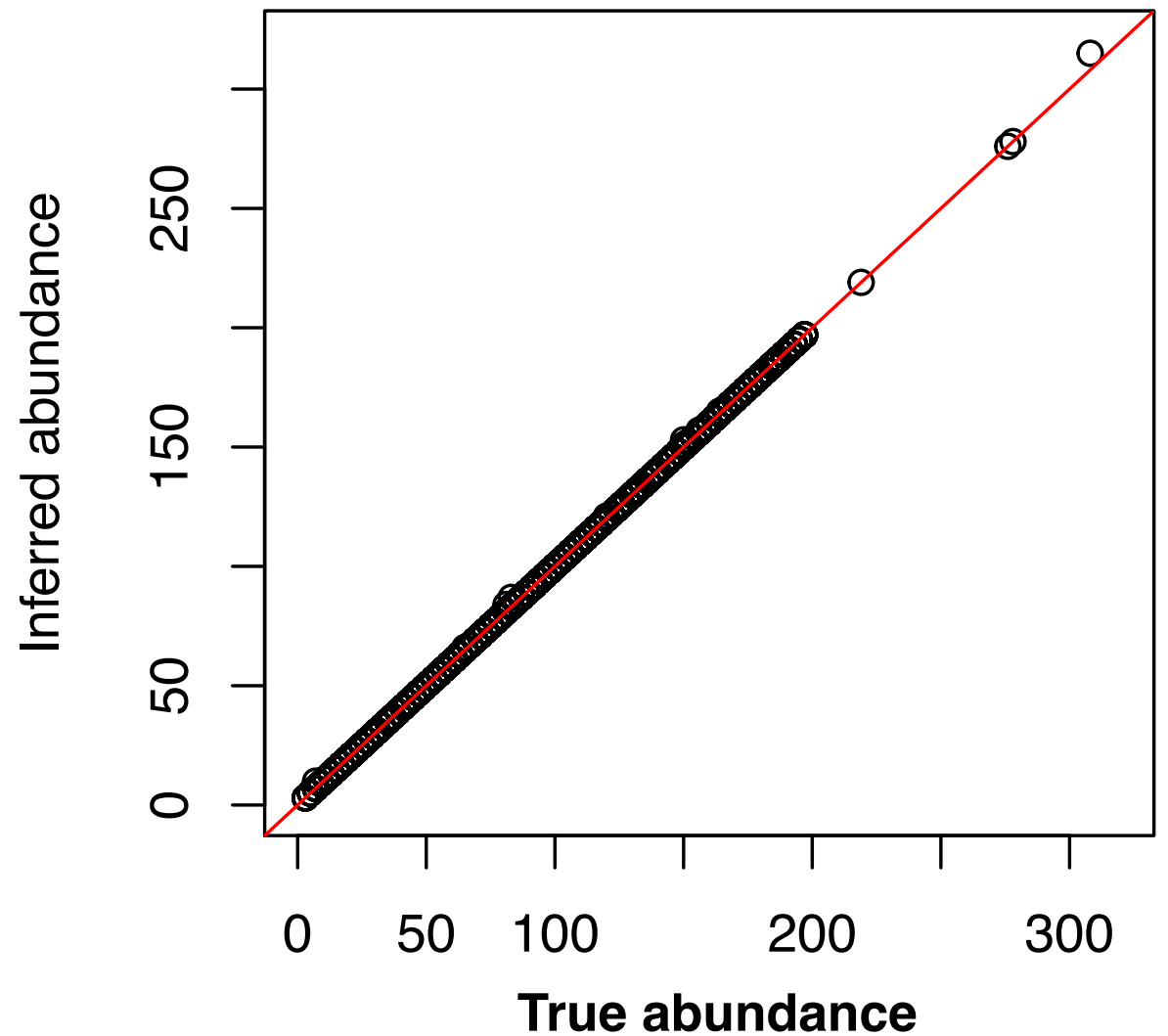
Accuracy: Simulated data

mothur (an)



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

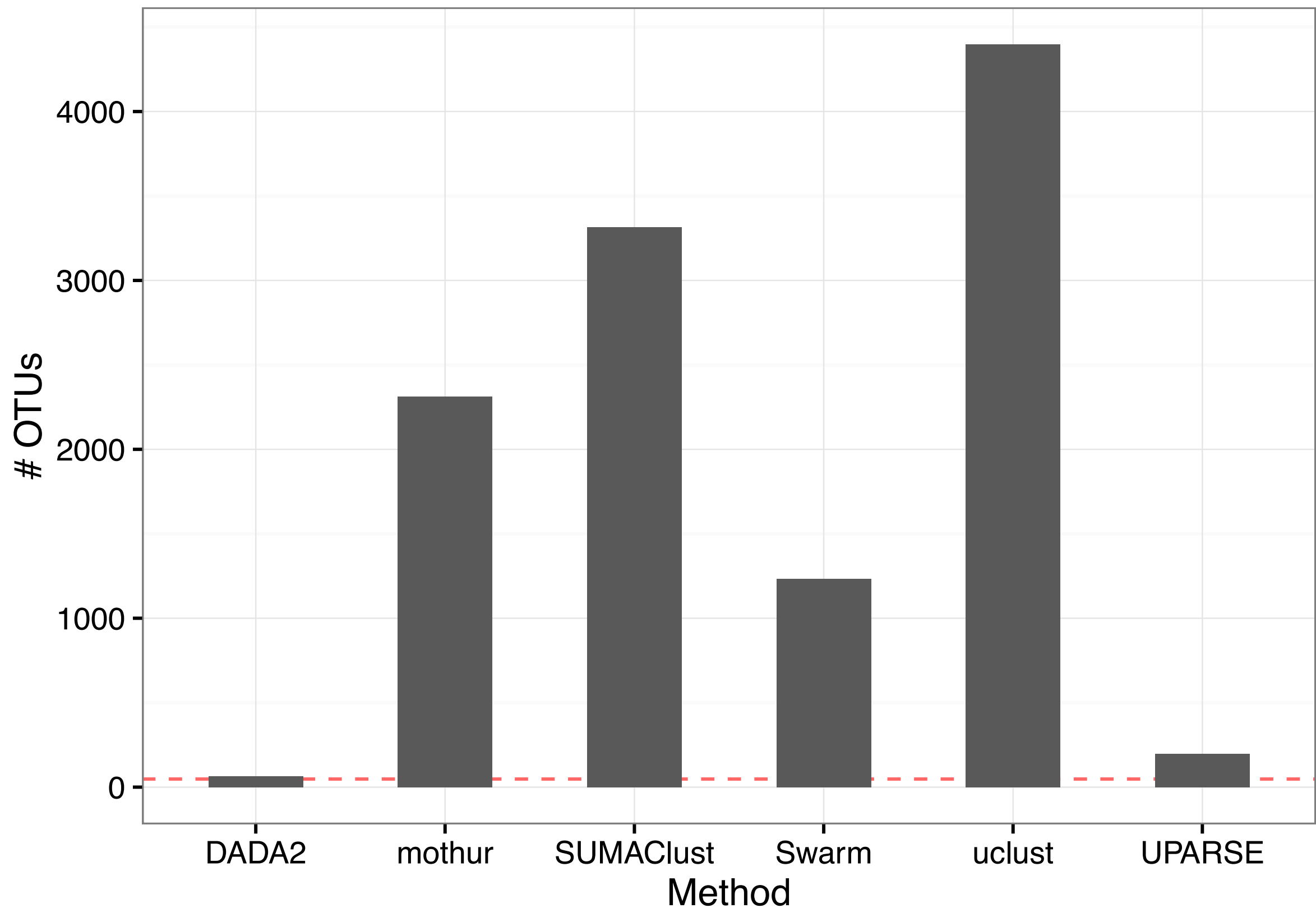
DADA2



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

Data: Kopylova, et al. mSystems, 2016.

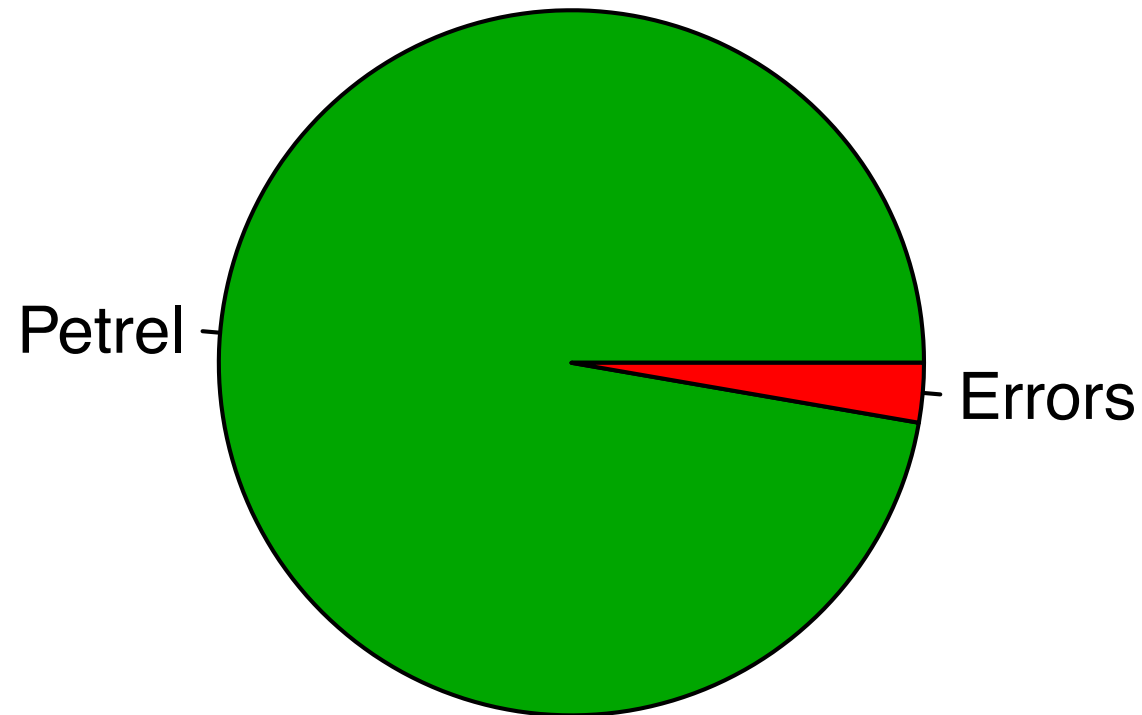
Accuracy: Mock community



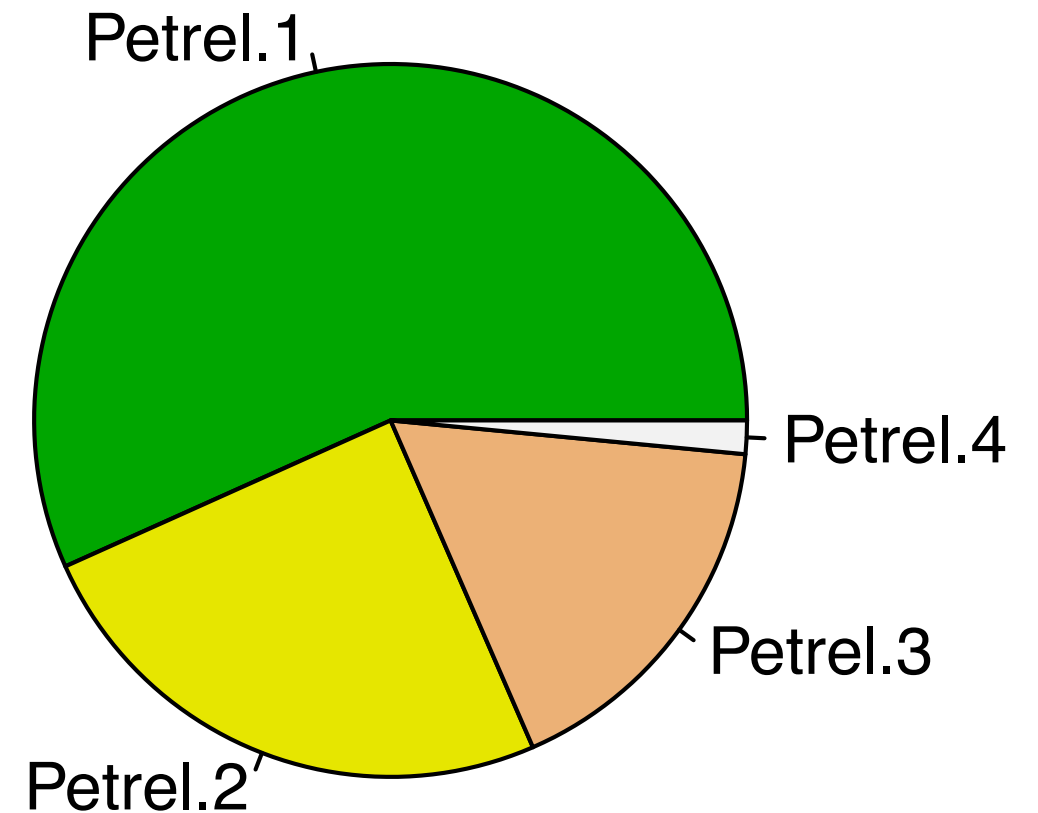
Credit: Kopylova, et al. mSystems, 2016.

Resolution: Petrel aDNA

QIIME: De novo



DADA2



Acknowledgements



Susan Holmes



Joey McMurdie



Michael Rosen



National Institutes of Health

<https://benjjneb.github.io/dada2/>