

# 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

NATURE METHODS | VOL.13 NO.7 | JULY 2016

Jeongsup Moon

## 1. Background

- A. Metagenome Analysis
- B. Limitations

## 2. Method

- A. Denoising
- B. Example

## 3. Results

- 1. Mock-up Community
- 2. Real Data

## 4. Conclusions

- 1. Limitation
- 2. Colab

# Metagenome analysis

FOOD

## We're Adding These Refreshing Probiotic Drinks to Our Routine ASAP—Here's Why

Dec. 7, 2021

SPONSORED



## Metagenome analysis

FOOD

We're Adding These Refreshing

P BEAUTY 

### A Why Postbiotics Work Better In Topical Skin Care Products Than Their Biotic Precursors



mbg Beauty Director

By Alexandra Engler 



# Metagenome analysis

FOOD

We're Addin

P BEAUTY

A Why Pos  
Care Pro

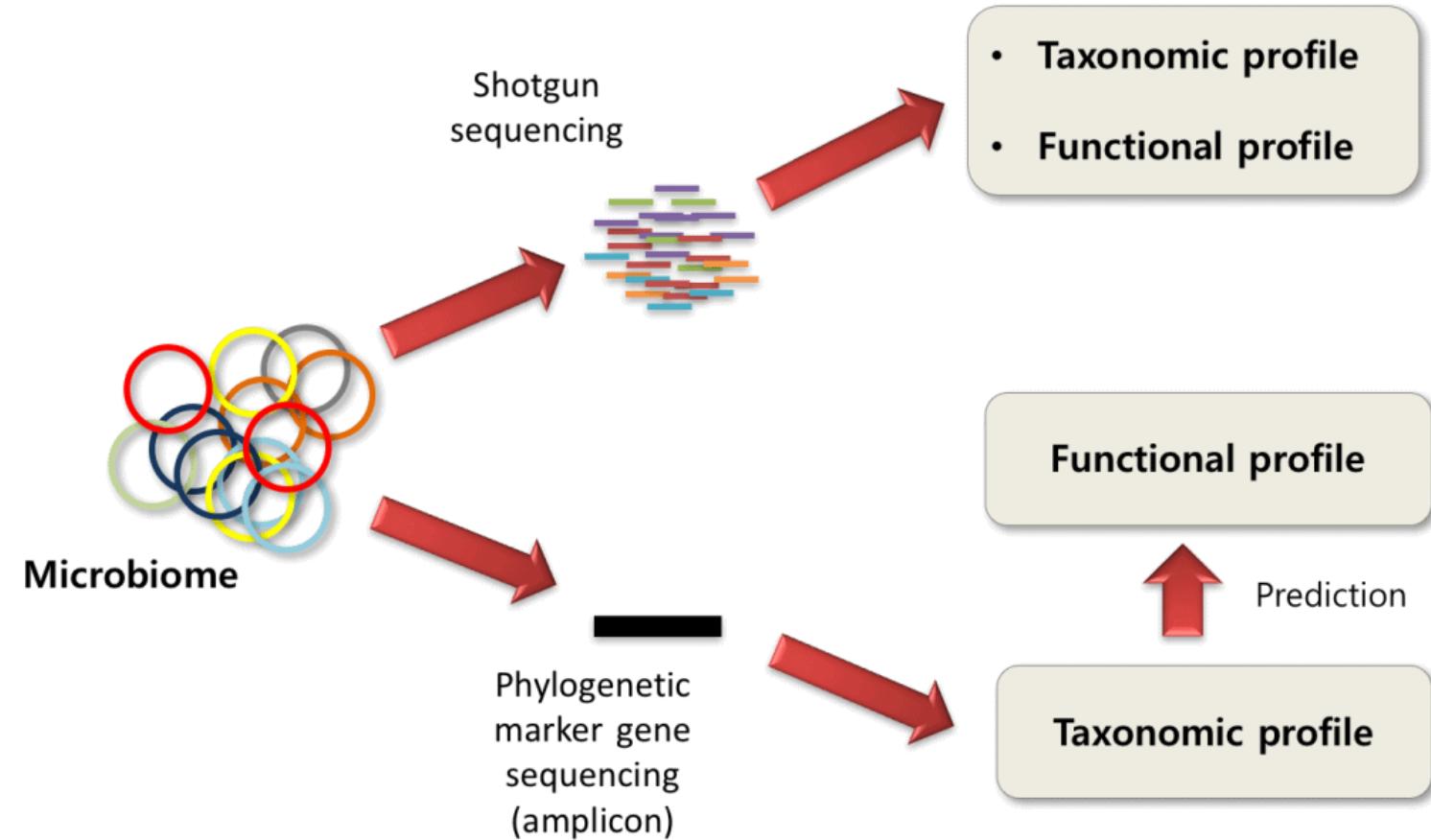
mbg Beauty  
By Alexandria

Boss Dog launches 11th product, probiotic-infused dog kibbles



Q. What is metagenome?

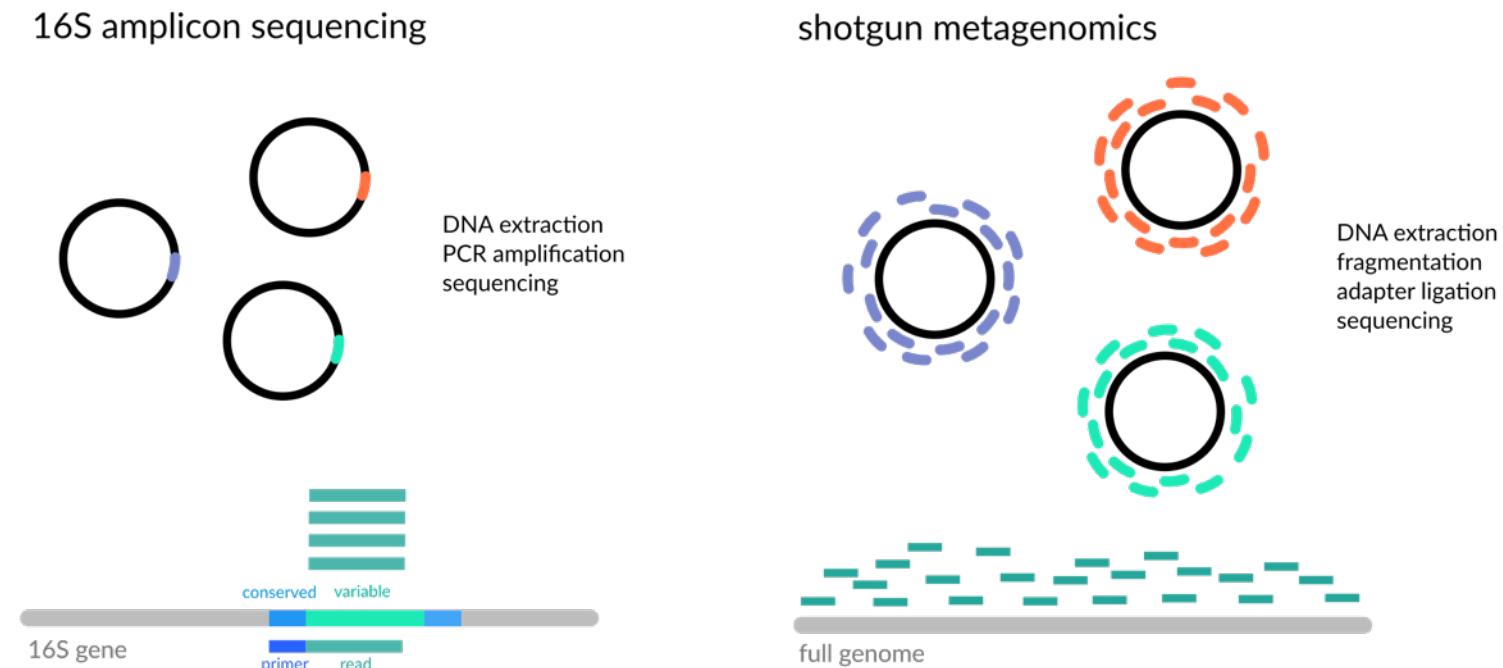
Q. How?



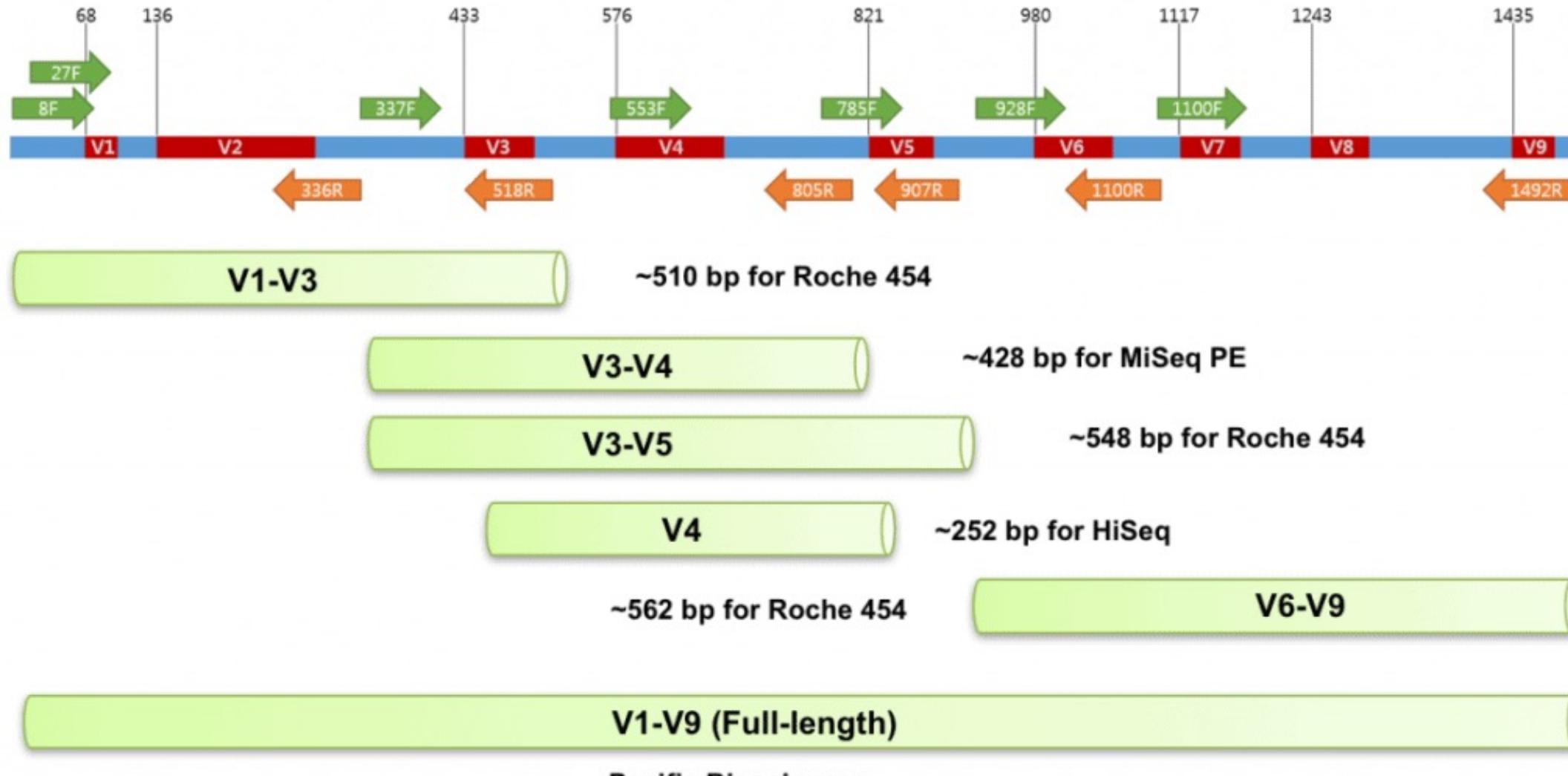
Q. What is metagenome?



Q. How?

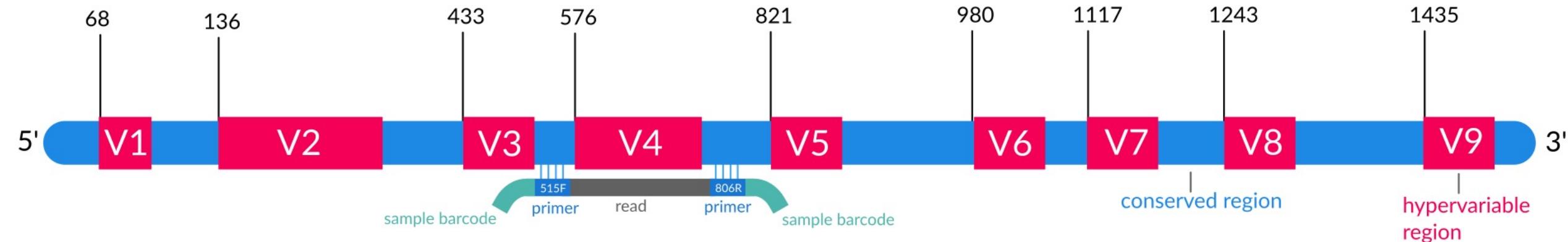


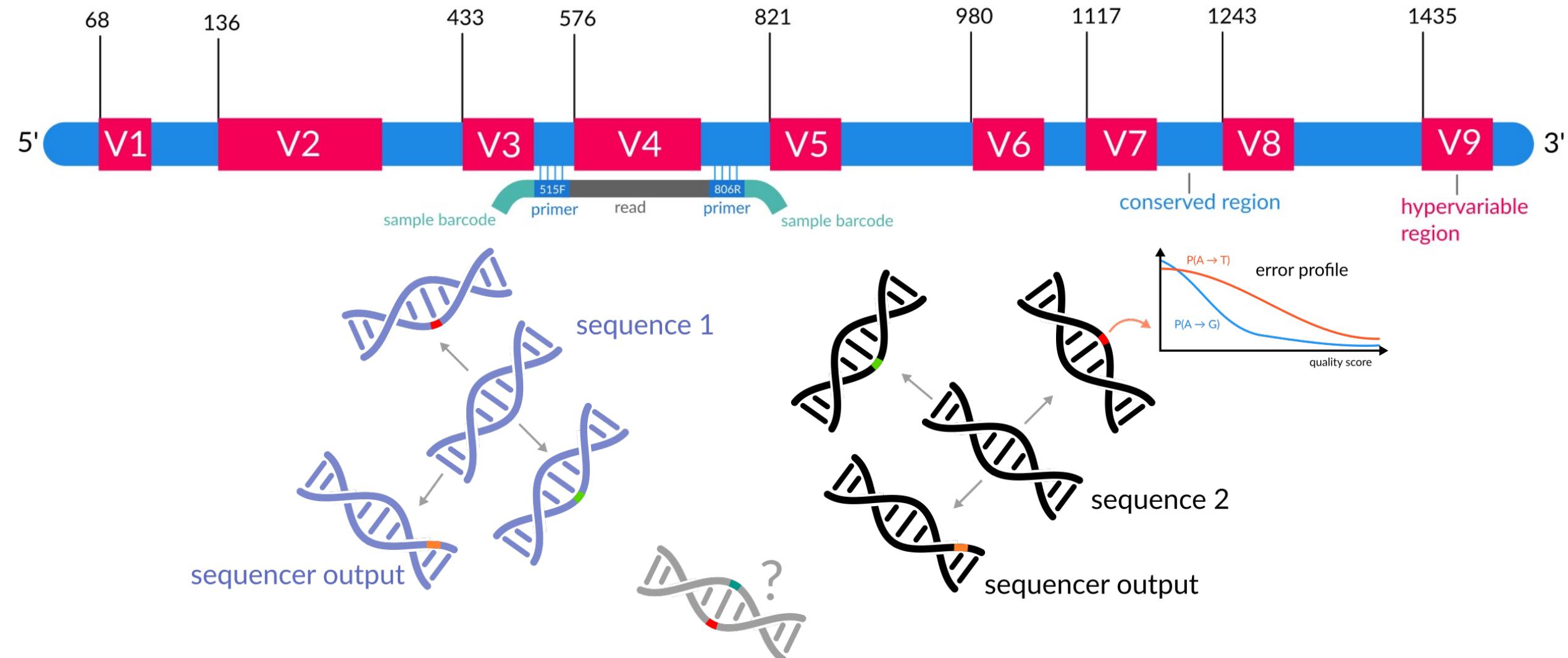
# Metagenome analysis

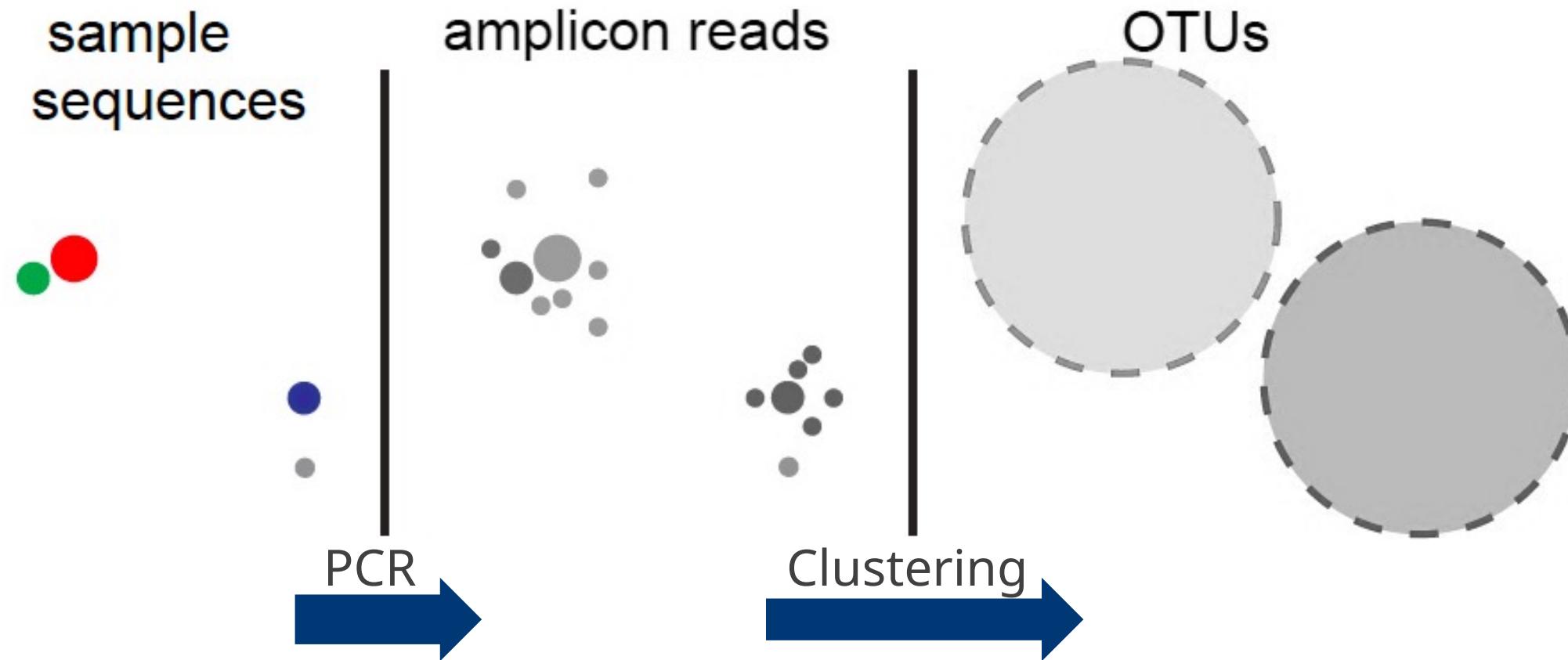


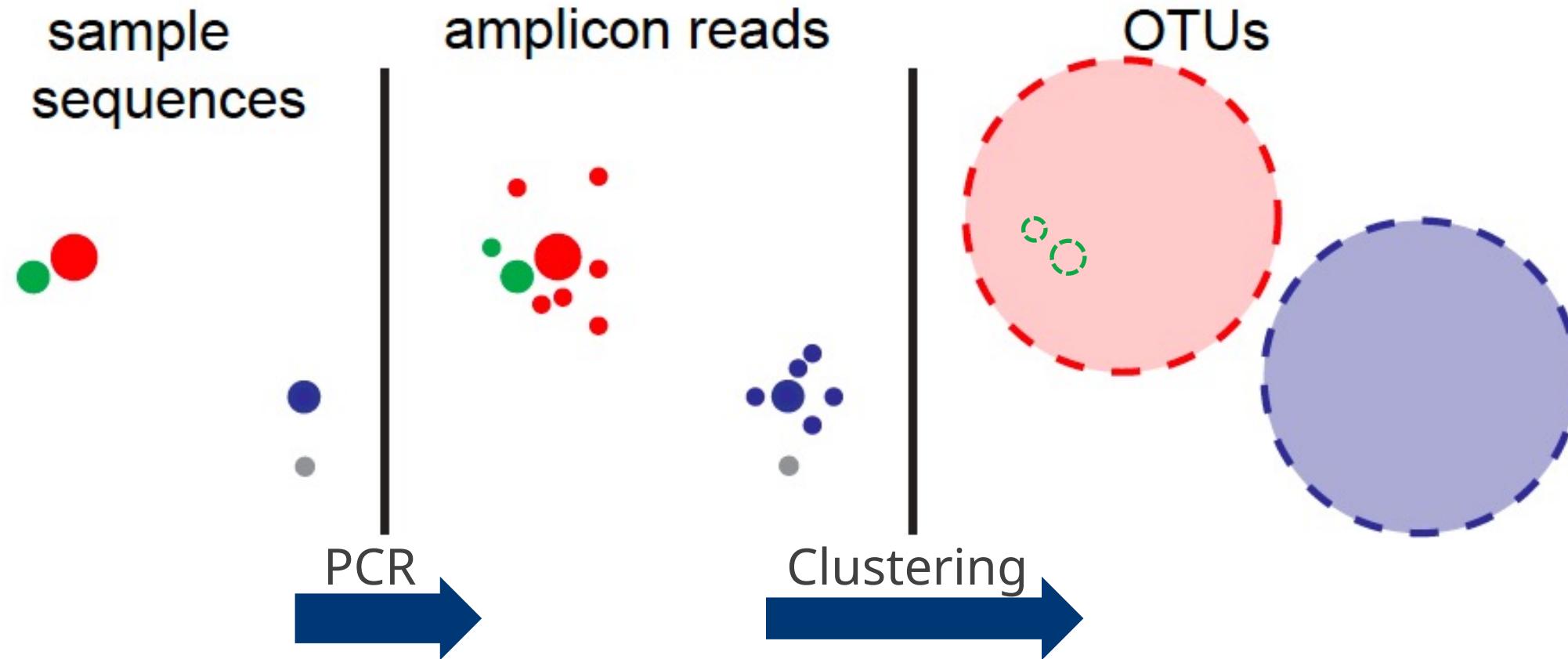
Pacific Biosciences

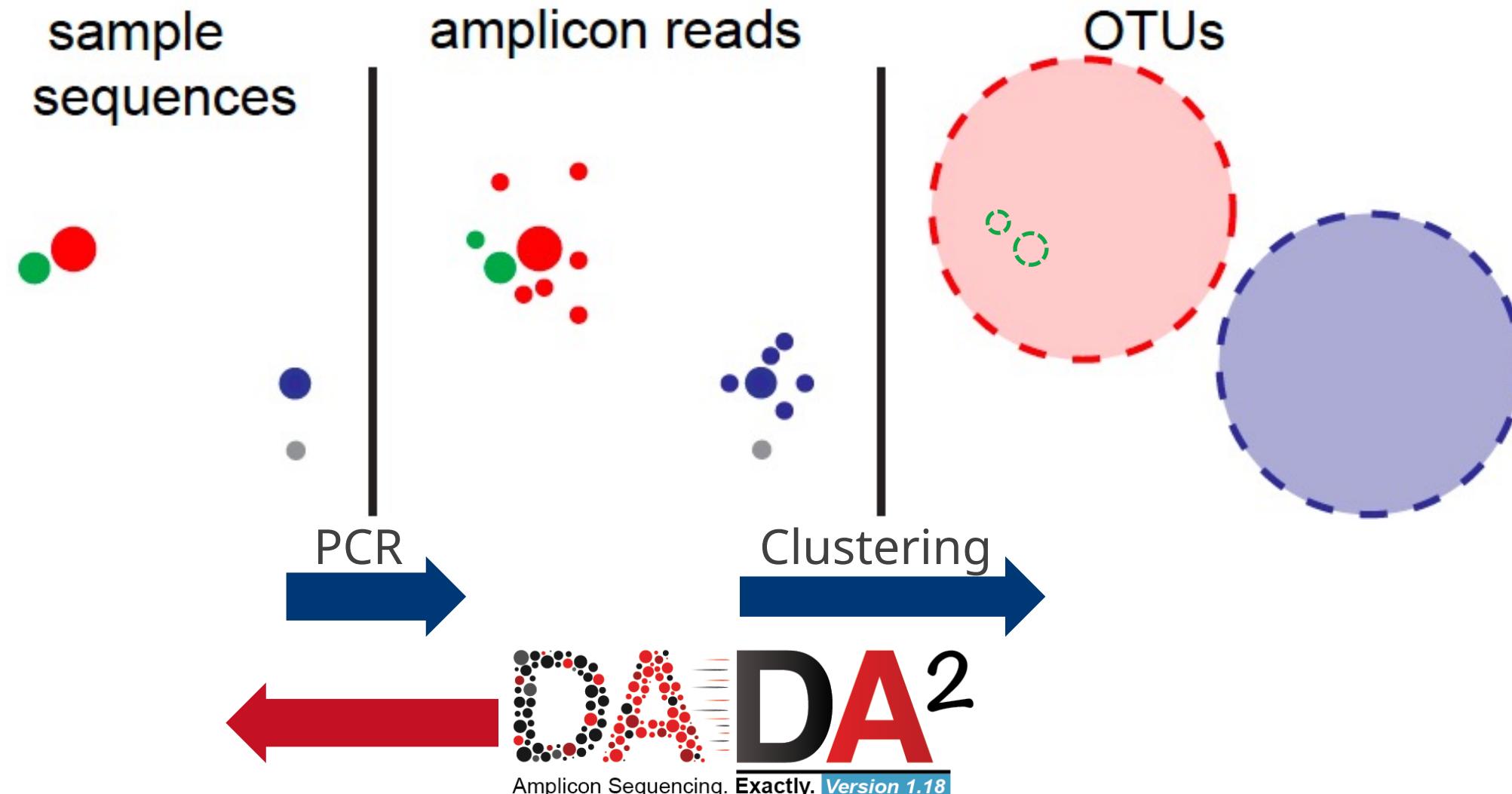
# Metagenome analysis











## DADA2: Two baseline assumptions

1. DADA2 models errors:
  - Independently within a read
  - Independantly between reads
  
2. DADA2 models abundance:
  - number of sequence  $i$  that is produced from sequence  $j$  will be Poisson distributed.

## DADA2: Denoising

DADA2 models errors:

- Independently within a read
- Independantly between reads

error rate

$$\lambda_{ji} = \prod_{l=0}^L p(j(l) \rightarrow i(l), \underline{qi(l)})$$

Probability of  $l$ th nt in sequence  $j$   
change to  $l$ th in sequence  $i$

Based on the quality

true sequence

$j$ : CCAACGGGAGACAGCAGTGGGGAA

error sequence

$i$ : CCTACGGGAG**G**CAGCAGTGGGGAA

Fastq File

- Seq\_id: @M00763:36:00000000-A8T0A:1:1101:1
- sequence : CCTACGGGAG**G**CAGCAGTGGGGAAATT
- delim: +
- quality: 88CCCGDBAF)====CEFFGGGG>GGGGGG

# DADA2: Denoising

DADA2 models errors:

- Independently within a read
- Independantly between reads

error rate

$$\lambda_{ji} = \prod_{l=0}^L p(j(l) \rightarrow i(l), q_i(l))$$

Probability of  $l$ th nt in sequence  $j$   
change to  $l$ th in sequence  $i$

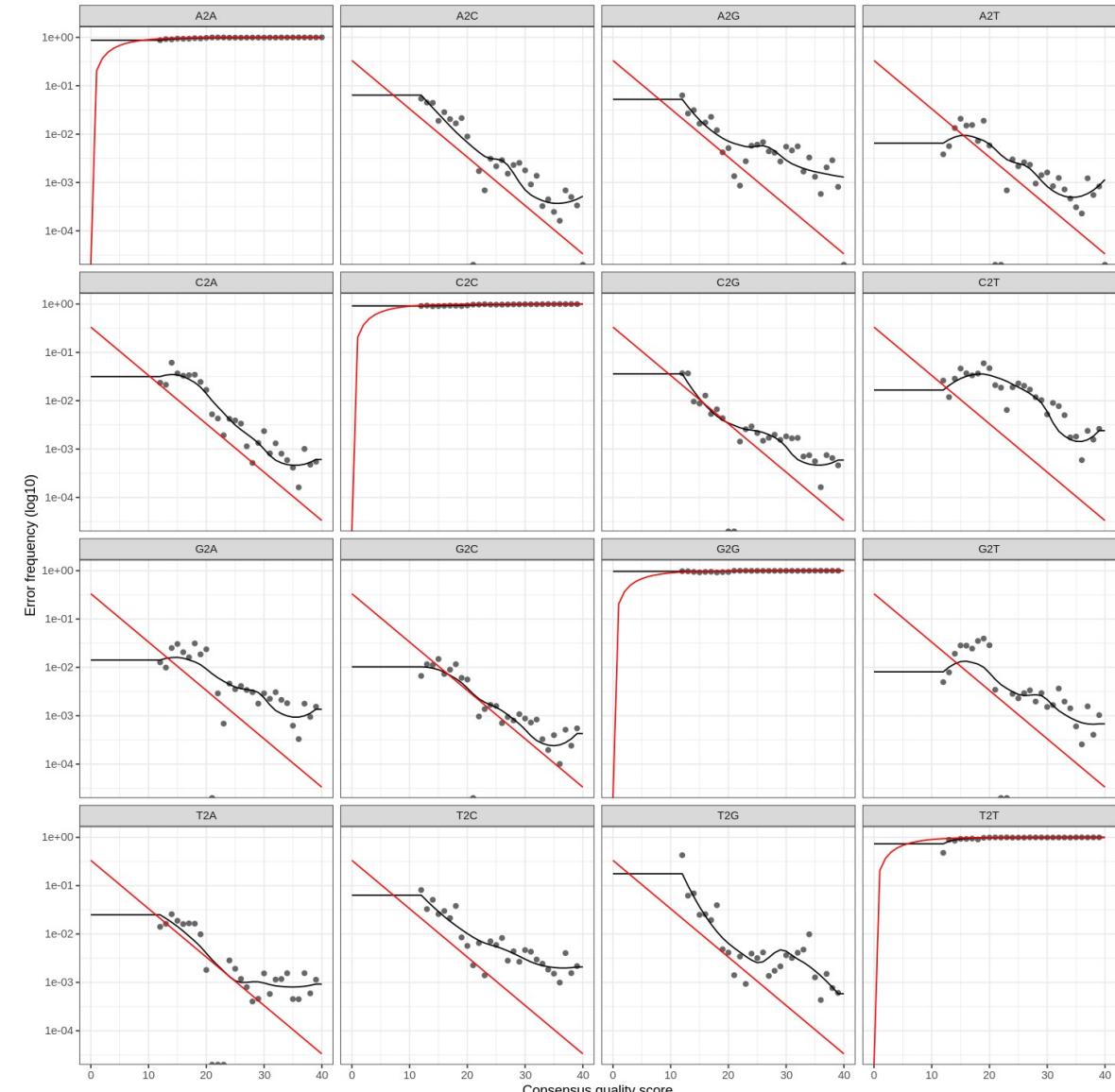
Based on the quality

true sequence

$j$ : CCAACGGGAGACAGCAGTGGGGAA

error sequence

$i$ : CCTACGGGAGG**C**AGCAGTGGGGAA



## DADA2: Denoising

DADA2 models abundance:

- number of sequence  $i$  that is produced from sequence  $j$  will be Poisson distributed

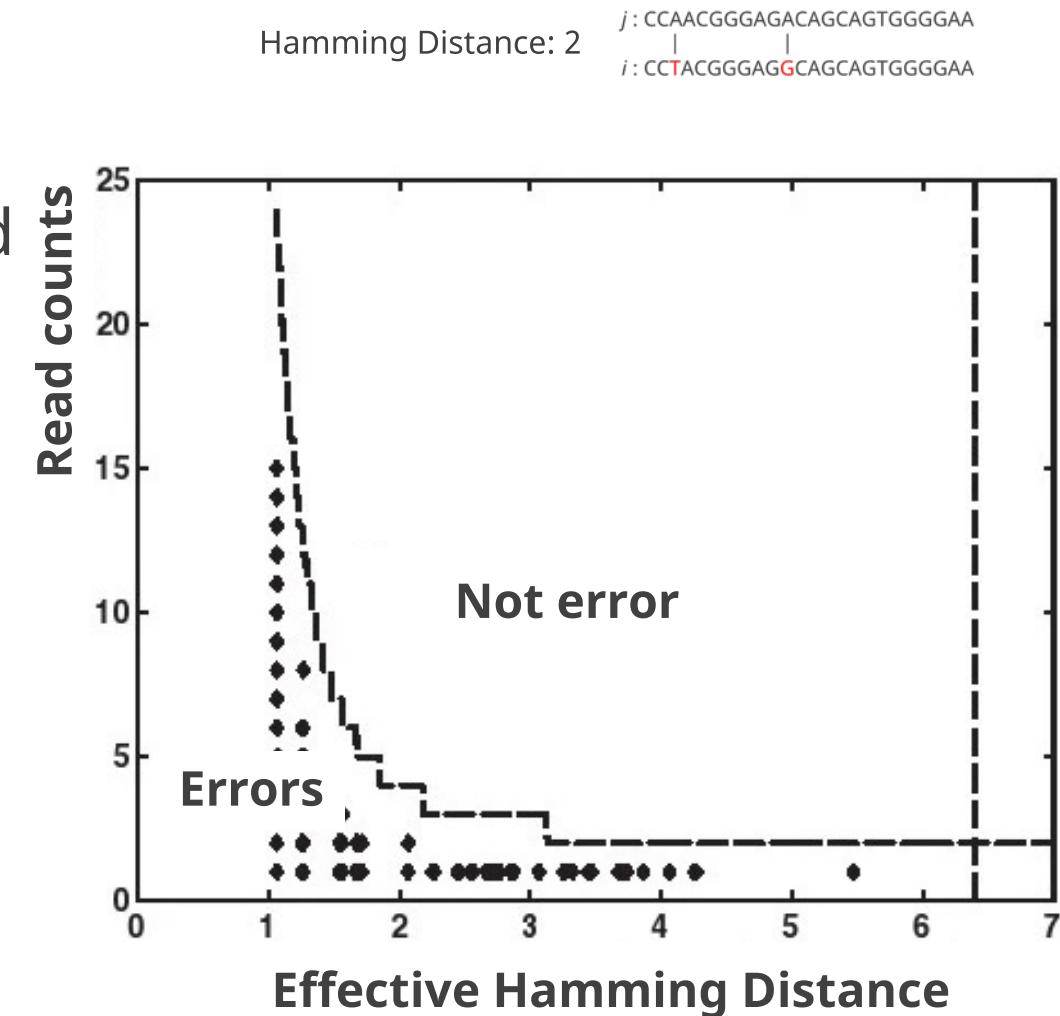
$$P(X = k) = \frac{\lambda^k e^{-\lambda}}{k!}$$

rate of events  
observe  $k$  events

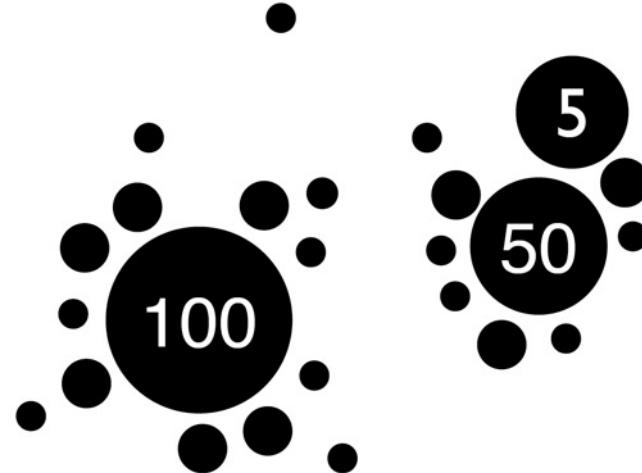
DADA2 models abundance probability:

$$p_A(j \rightarrow i) = \frac{\sum_{a=a_i}^{\infty} \rho_{pois}(n_j \lambda_{ji}, a)}{1 - \rho_{pois}(n_j \lambda_{ji}, 0)}$$

# of sequence  $i$   
error rate  
# of sequence  $j$

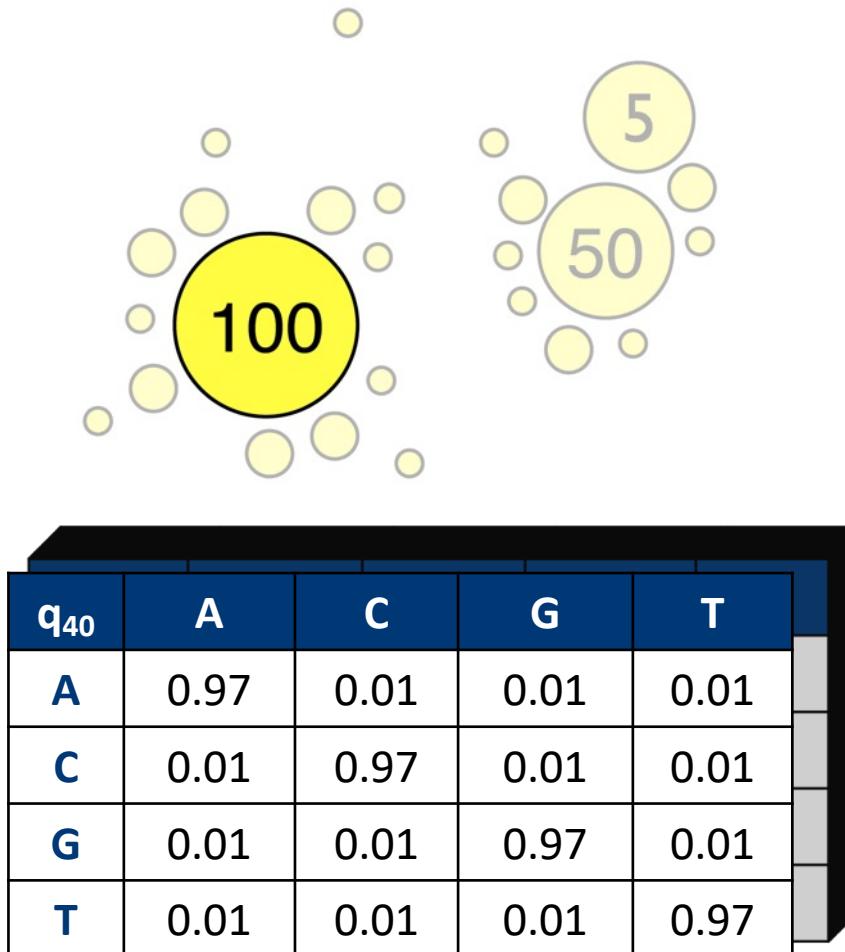


## DADA2: Example



1. Assumes : The most abundant sequence as **the only real sequence  $j$** , and others **as errors  $i$** .

## DADA2: Example

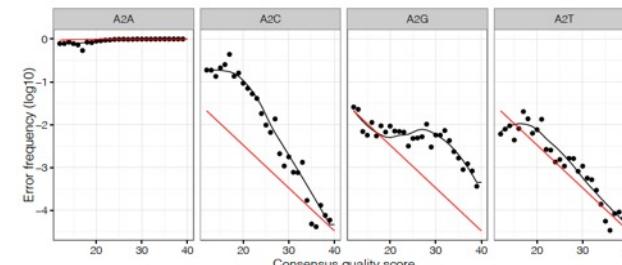


\* This is a simplified version of 16 \* 40 table.

\* Note that there are 40 different q values.

1. Assumes : The most abundant sequence as **the only real sequence  $j$** , and others as **errors  $i$** .

2. Learns error model parameter(transition probability).  
Learns error model, and calculate error rate.  
If mismatch > 10%,  $\lambda_{ji}$  is set to 0.

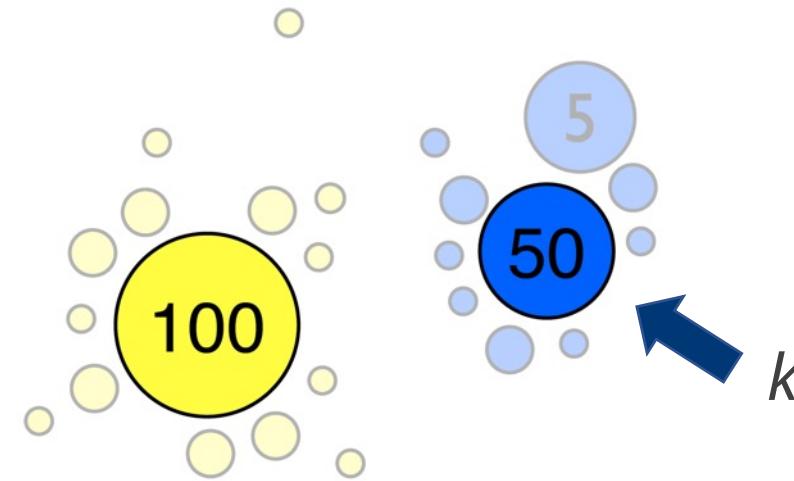


$$\lambda_{ji} = \prod_{l=0}^L p(j(l) \rightarrow i(l), q_i(l))$$

3. Calculates abundance p value.

$$p_A(j \rightarrow i) = \frac{\sum_{a=a_i}^{\infty} \rho_{pois}(n_j \lambda_{ji}, a)}{1 - \rho_{pois}(n_j \lambda_{ji}, 0)}$$

## DADA2: Example



<b>q<sub>40</sub></b>	<b>A</b>	<b>C</b>	<b>G</b>	<b>T</b>
<b>A</b>	0.997	0.001	0.001	0.001
<b>C</b>	0.001	0.997	0.001	0.001
<b>G</b>	0.001	0.001	0.997	0.001
<b>T</b>	0.001	0.001	0.001	0.997

\* This is a simplified version of 16 \* 40 table.

\* Note that there are 40 different q values.

3. Calculates abundance p value.

$$p_A(j \rightarrow i) = \frac{\sum_{a=1}^{\infty} a_i \rho_{pois}(n_j \lambda_{ji}, a)}{1 - \rho_{pois}(n_j \lambda_{ji}, 0)}$$

4. Select lowest abundance p value, and set it as another real **sequence k**.

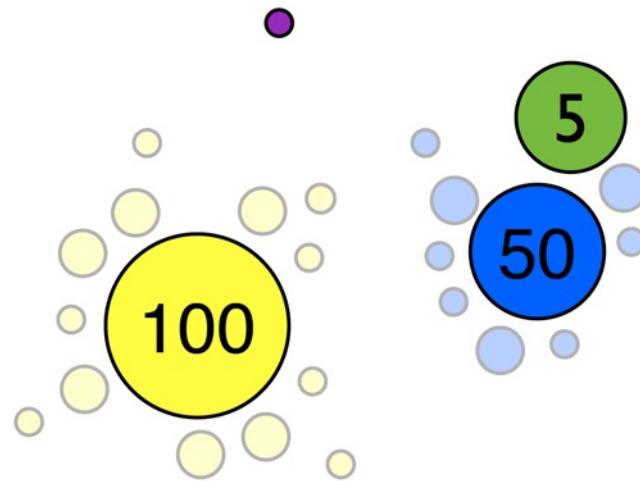
5. Align errors *i* to where they will have most expected number.  $E[\rho_{pois}(n_j; \lambda_{ji})]$  or  $E[\rho_{pois}(n_k; \lambda_{ki})]$

6. Update error models

$$\lambda_{ji} = \prod_{l=0}^L p(j(l) \rightarrow i(l), q_i(l))$$

$$\lambda_{ki} = \prod_{l=0}^L p(k(l) \rightarrow i(l), q_i(l))$$

## DADA2: Example



<b>q<sub>40</sub></b>	<b>A</b>	<b>C</b>	<b>G</b>	<b>T</b>
<b>A</b>	0.998	1x10 <sup>-4</sup>	2x10 <sup>-3</sup>	2x10 <sup>-4</sup>
<b>C</b>	6x10 <sup>-5</sup>	0.999	3x10 <sup>-6</sup>	1x10 <sup>-3</sup>
<b>G</b>	1x10 <sup>-3</sup>	3x10 <sup>-6</sup>	0.999	6x10 <sup>-5</sup>
<b>T</b>	2x10 <sup>-4</sup>	6x10 <sup>-3</sup>	1x10 <sup>-4</sup>	0.998

\* This is a simplified version of 16 \* 40 table.

\* Note that there are 40 different q values.

3. Calculates abundance p value.

$$p_A(j \rightarrow i) = \frac{\sum_{a=1}^{\infty} a_i \rho_{pois}(n_j \lambda_{ji}, a)}{1 - \rho_{pois}(n_j \lambda_{ji}, 0)}$$

4. Select lowest abundance p value, and set it as another real **sequence k**.

5. Align errors *i* to where they will have most expected number.  $E[\rho_{pois}(n_j; \lambda_{ji})]$  or  $E[\rho_{pois}(n_k; \lambda_{ki})]$

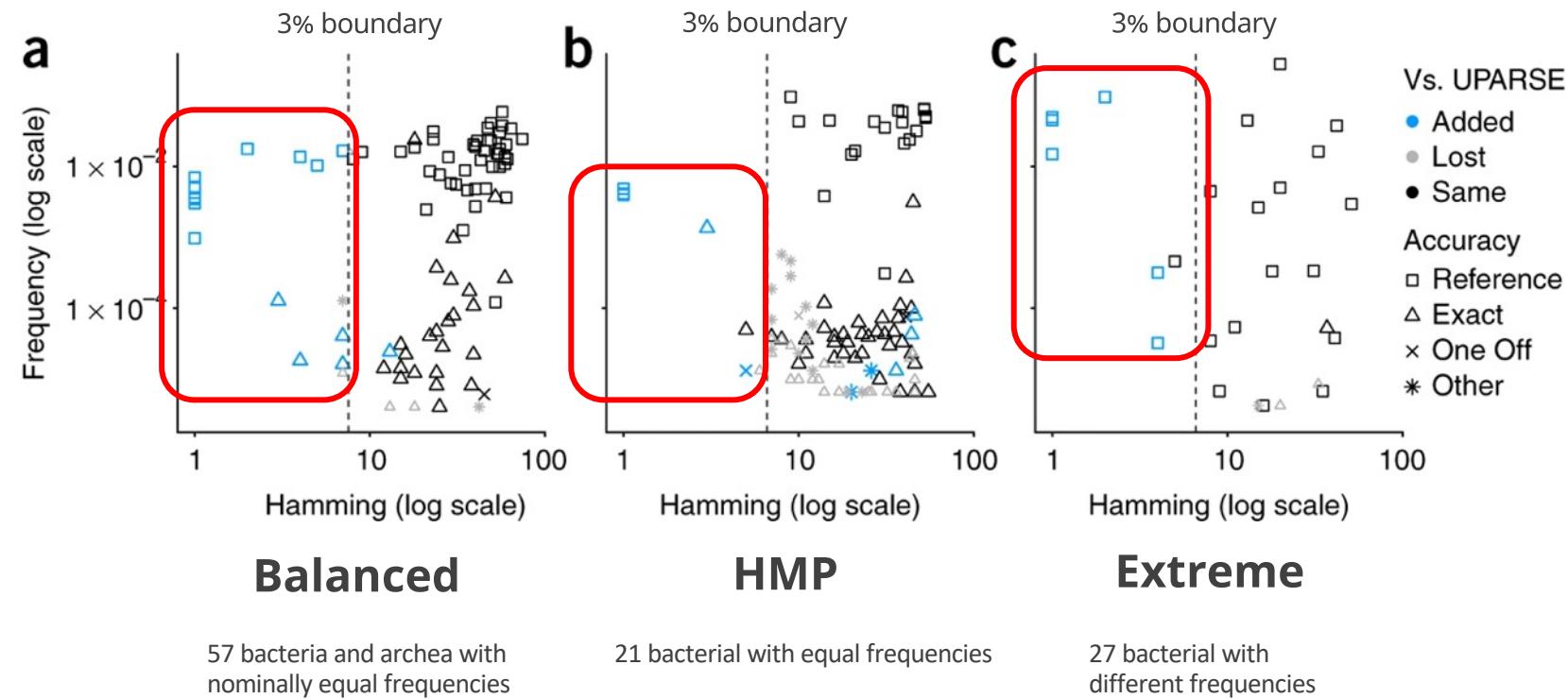
6. Update error models

$$\lambda_{ji} = \prod_{l=0}^L p(j(l) \rightarrow i(l), q_i(l))$$

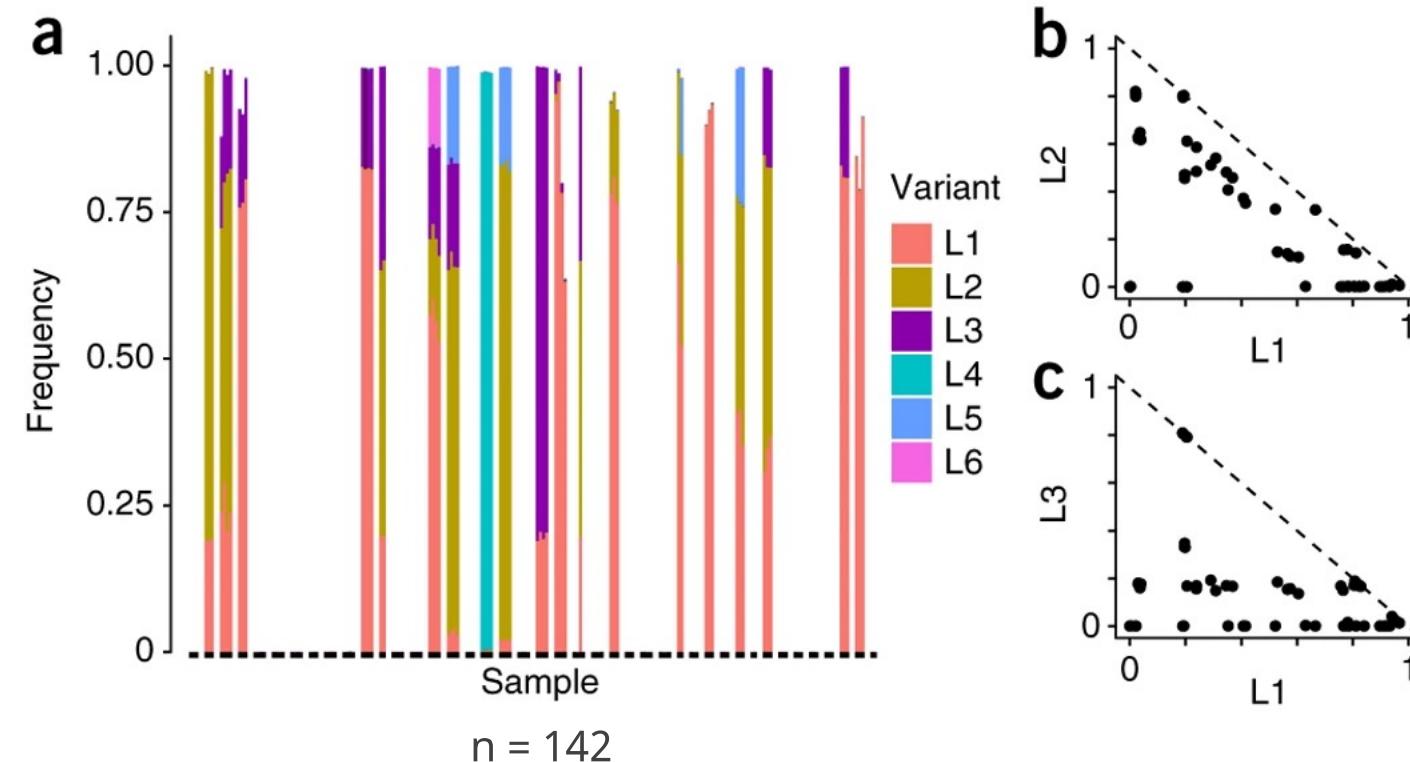
7. Do step 3 to 6, until all abundance p-value are above threshold( $1e^{-40}$ ).

# Mock-up Community

## Comparison of sequence variants inferred by DADA2 with OTUs constructed by UPARSE.



### Frequency of *Lactobacillus crispatus* (Strain level)



# Limitation

Strains missed by DADA2	Total Reads	Max Abundance	Hamming to NN	NN Abundance
<b>Prevotella buccalis</b>	5	1	51	9
<b>Clostridium methylpentosum DSM 5476</b>	5	1	25	13
<b>Clostridium phytofermentans ISDg</b>	11	2	15	59564
<b>Parabacteroides sp. D13</b>	3	2	1	28242

$$p_A(j \rightarrow i) = \frac{\sum_{a=a_i}^{\infty} \rho_{pois}(n_j \lambda_{ji}, a)}{1 - \rho_{pois}(n_j \lambda_{ji}, 0)}$$

Too many # of sequence  $j$

Singleton removed

[https://colab.research.google.com/drive/1CaLLzB1gEEuMk4P8J\\_siDFgLUX1d3gor?usp=sharing](https://colab.research.google.com/drive/1CaLLzB1gEEuMk4P8J_siDFgLUX1d3gor?usp=sharing)

# Thank you