

# DNA sequence analysis: data handling, visualization and (some) multivariate statistics in R

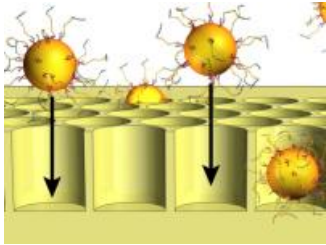
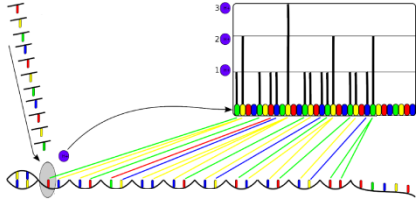
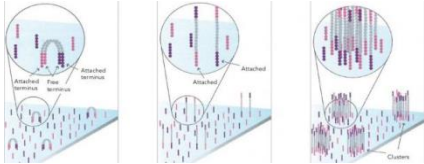
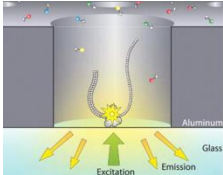
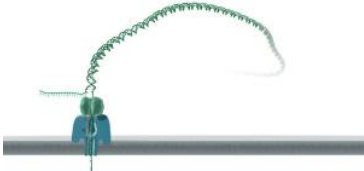
R roundtable 21.1.2016 – Christiane Hassenrück

# outline

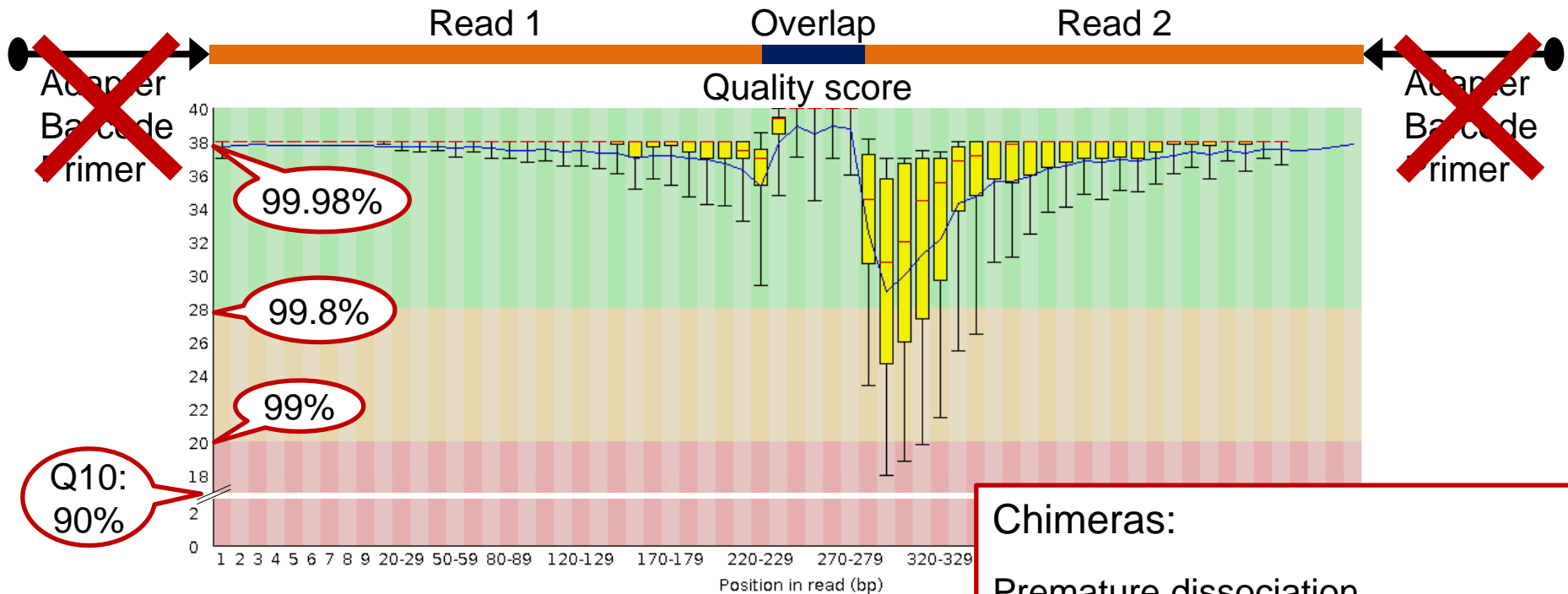
---

- |  |  |
|--|--|
| 1. Sequencing technologies and sequence analysis pipelines | Major steps, sequencing errors, OTU clustering algorithms, taxonomy  |
| 2. Data handling in R                                      | Reading data: Data subsetting, parsing taxonomic paths   |
| 3. Diversity concepts                                      | Alpha and beta diversity, random subsampling, effective species number, symmetric vs. asymmetric diversity indices |
| 4. Plotting community data in R                            | Custom functions (alpha diversity, abundant taxa, networks), ordination plots                                      |
| 5. Pitfalls of sequence analysis                           | Compositionality, sampling design and replication  |
| 6. Multivariate statistics                                 | Overall patterns (ANOSIM, PERMANOVA, RDA), differential OTU abundance  |
| 7. Estimating bacterial functions                          | Tax4Fun  |
-

# NGS technologies

Sequencing technology	Principle	Read length	Errors	Comments
454		< 450 bp SE	homopolymers	discontinued
Ion torrent		< 400 bp SE	homopolymers	
<b>Illumina</b>		< 300 bp SE < 550 bp PE	substitutions	Currently preferred for amplicon and shotgun sequencing
PacBio		> 10 kb	~ 10% error rate (single pass)	
nanopore		> 10 kb	< 4% error rate	In development

# Bioinformatic analysis



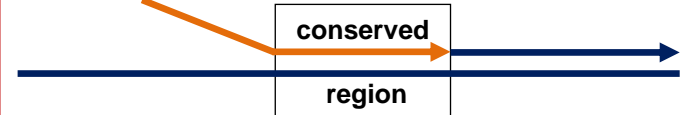
- 1) Adapter, barcode, primer removal
- 2) Quality trimming
- 3) Merging
- 4) OTU clustering
- 5) Taxonomic classification

## Chimeras:

Premature dissociation



Annealing to new template



Chimeric sequence



# OTU clustering methods

Algorithm	Pro's	Con's
Hierarchical	<ul style="list-style-type: none"><li>• Better defined OTUs than heuristic clustering</li></ul>	<ul style="list-style-type: none"><li>• Very slow</li></ul>
Heuristic (greedy)	<ul style="list-style-type: none"><li>• Fast compared to hierarchical clustering</li></ul>	<ul style="list-style-type: none"><li>• Low reproducibility</li></ul>
Swarm	<ul style="list-style-type: none"><li>• Fast</li><li>• Variable OTU cut-off</li><li>• High reproducibility</li></ul>	<ul style="list-style-type: none"><li>• Large swarms</li></ul>
Oligotype	<ul style="list-style-type: none"><li>• Fast</li><li>• Omits stochastic variation</li><li>• Sub-species resolution (SNPs)</li></ul>	<ul style="list-style-type: none"><li>• No rare biosphere</li></ul>

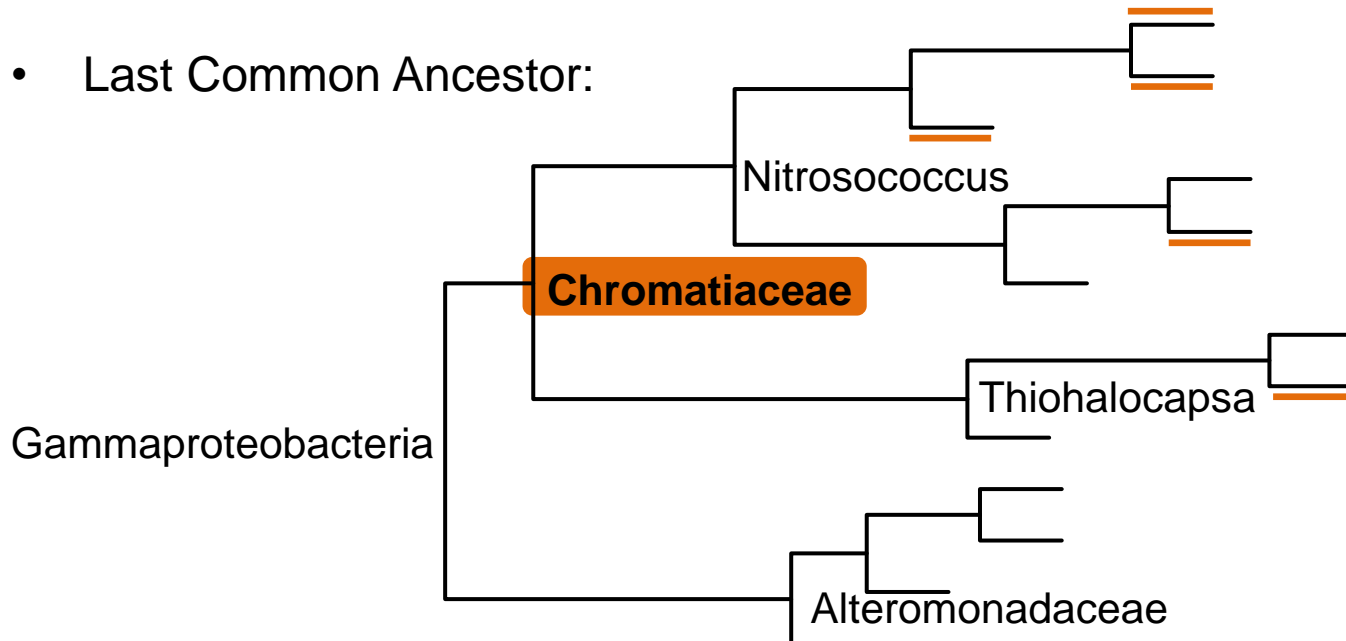
# Taxonomic classification

Domain;Phylum;Class;Order;Family;Genus

←   ←   ←   ←   ←

Bacteria;Proteobacteria;Gammaproteobacteria;Chromatiales;Chromatiaceae;Nitrosococcus

- Last Common Ancestor:



- Truncated paths:  
Bacteria;Proteobacteria;Gammaproteobacteria;Chromatiales;Chromatiaceae;**unclassified**
- Incomplete paths:  
Bacteria;Proteobacteria;Gammaproteobacteria;**Incertae Sedis**;**Incertae Sedis**;Sedimenticola

Let's get to 

Sample-by-OTU table  
OTU-by-taxonomy table  
Environmental data

**Consistent  
sample and OTU  
names!**

# Alpha diversity


- Community richness and evenness per sample
- Classical indices:
  - OTU number
  - Shannon
  - Chao 1
  - Inverse Simpson
- Hill numbers: effective species number
  - 1 formula for all indices → only changing 1 parameter (q)

$${}^qD = \left( \sum_{i=1}^S p_i^q \right)^{1/(1-q)}$$

**${}^0D$  = OTU number**

**${}^1D$  = exp (Shannon)**

**${}^2D$  = inverse Simpson**



Influence of  
rare biosphere

- Rarefaction curves
- Unequal sequencing depth → random subsampling to compare alpha diversity indices across samples



# Beta diversity

- Community (dis)similarity between samples

	OTU1	OTU2	OTU3	OTU4		OTU1	OTU2	OTU3	OTU4	
S1	14	2	14	14	presence/ absence	S1	1	1	1	1
S2	10	14	0	8		S2	1	1	0	1
S3	0	5	0	2		S3	0	1	0	1
S4	0	0	1	0		S4	0	0	1	0

Asymmetrical vs. symmetrical  
Bray-Curtis vs. euclidean

Jaccard

	S1	S2	S3	S4		S1	S2	S3	S4		S1	S2	S3	S4
S1	0				S1	0				S1	0			
S2	0.5	0			S2	19.8	0			S2	0.25	0		
S3	0.8	0.6	0		S3	23.3	14.7	0		S3	0.5	0.33	0	
S4	1.0	1	1	0	S4	23.8	19	5.5	0	S4	0.75	1	1	0

- Zeros in ecology: Is this species really not there or did we just not find it?  
→ double zeros not relevant

# Plotting in

Alpha diversity indices

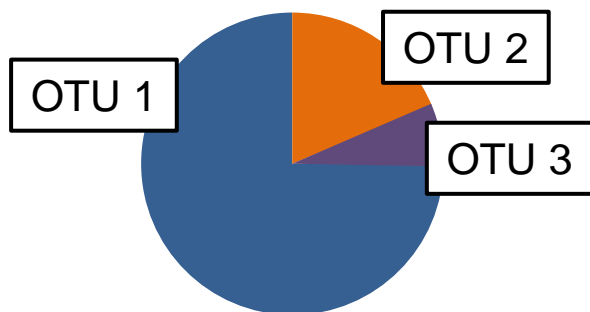
Abundant taxa

OTU networks

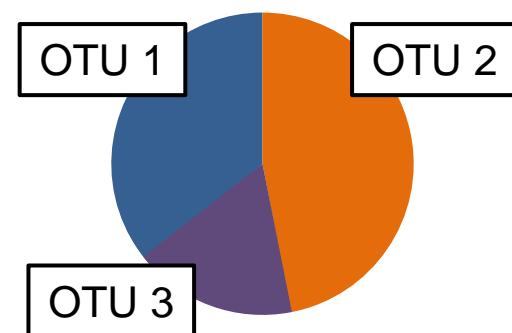
Ordination plots

# Pitfalls of sequence analysis

- Compositionality: OTU abundances not independent
- Centered log-ratio transformation (clr):  $\log(x_i) - \log(n\sqrt{\text{product}(x_1 \dots x_n)})$



Decrease in  
OTU 1



Sample 1	%	clr
OTU 1	75	1.83
OTU 2	18	-0.23
OTU 3	7	-1.59

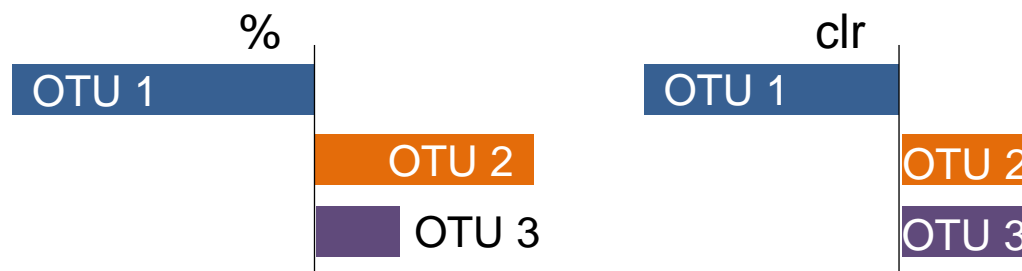
%:  $\Delta 11$   
 clr:  $\Delta 1.36$

Sample 2	%	clr
OTU 1	35	0.14
OTU 2	47	0.61
OTU 3	18	-0.75

%:  $\Delta 29$   
 clr:  $\Delta 1.36$

→ Difference between OTUs independent of library size

- Difference between samples:



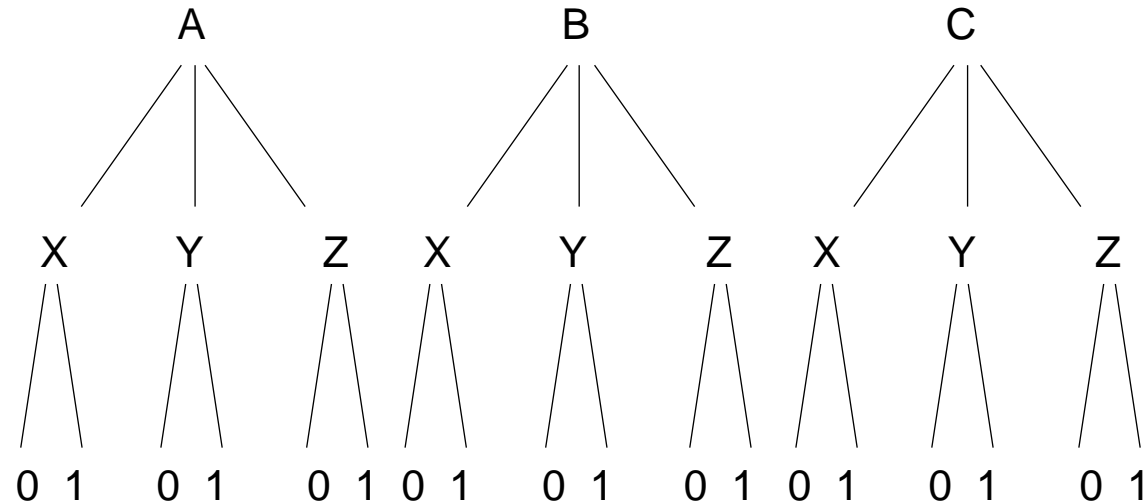
# Pitfalls of sequence analysis

- Sampling design

Factor 1:

Factor 2:

Factor 3:



Replication:

Estimated costs:

$$3 \times 3 \times 2 \times 5 = 90$$

~ 5000 € consumables

- To reduce costs, i.e. only sequencing selected samples:

**Better to remove a condition than to reduce the number of replicates!**

# Multivariate statistics in



ANOSIM and PERMANOVA

Redundancy analysis

Differential OTU abundance (ALDEx2)

Path analysis

→ check out: <http://mb3is.megx.net/gustame>

# Estimating bacterial functions



Tax4Fun: <http://tax4fun.gobics.de/>

Sample-by-OTU table

Taxonomic paths

Taxonomic reference database

# Useful links

- Web links:
  - <http://www.arb-silva.de/>
  - <http://www.arb-silva.de/download/archive/>
  - <http://mb3is.megx.net/gustame>
  - <http://tax4fun.gobics.de/>
  - <https://github.com/chassenr/NGS>
- References:
  - Aßhauer, K.P., Wemheuer, B., Daniel, R., and Meinicke, P. (2015) Tax4Fun : predicting functional profiles from metagenomic 16S rRNA data. *Bioinformatics* **31**: 2882–2884.
  - Chao, A., Gotelli, N.J., Hsieh, T.C., Sander, E.L., Ma, K.H., Colwell, R.K., and Ellison, A.M. (2014) Rarefaction and extrapolation with Hill numbers: A framework for sampling and estimation in species diversity studies. *Ecol. Monogr.* **84**: 45–67.
  - Eren, A.M., Morrison, H.G., Lescault, P.J., Reveillaud, J., Vineis, J.H., and Sogin, M.L. (2014) Minimum entropy decomposition : Unsupervised oligotyping for sensitive partitioning of high- throughput marker gene sequences. *Isme J* **9**: 968–979.
  - Fernandes, A.D., Reid, J.N., Macklaim, J.M., McMurrough, T.A., Edgell, D.R., and Gloor, G.B. (2014) Unifying the analysis of high-throughput sequencing datasets: characterizing RNA-seq, 16S rRNA gene sequencing and selective growth experiments by compositional data analysis. *Microbiome* **2**: 15.
  - Mahé, F., Rognes, T.T., Quince, C., de Vargas, C., and Dunthorn, M. (2014) Swarm : robust and fast clustering method for amplicon-based studies. *PeerJ* 1–12.
  - Sinclair, L., Osman, O.A., Bertilsson, S., Eiler, A. (2015) Microbial Community Composition and Diversity via 16S rRNA Gene Amplicons: Evaluating the Illumina Platform. *PLoS ONE* **10**: e0116955.