

BIO9905MERG1 – Bioinformatics for Environmental Sequencing (DNA metabarcoding)

Welcome!

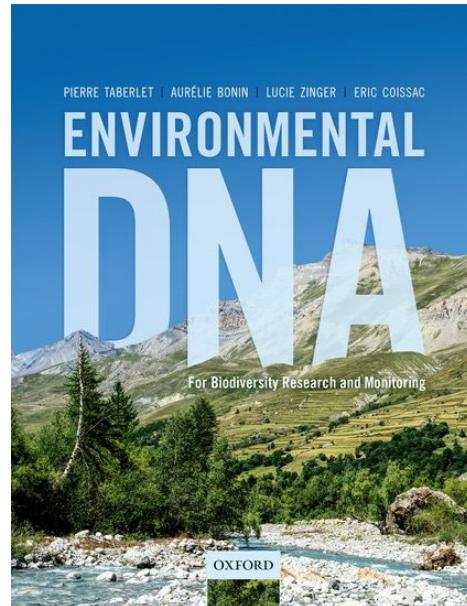


Introduction to DNA metabarcoding

- Explain terms
- Introduce key steps
- Introduce some literature
- More in-depth information in later talks



Pierre Taberlet



Some important / confusing terms

ASVs	amplicon sequencing
DNA metabarcoding	OTUs
eDNA	metatranscriptomics
metagenetics	metagenomics
MOTUs	environmental sequencing
marker gene analysis	microbiome analysis
	community profiling

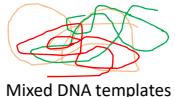
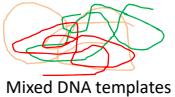
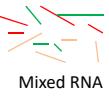
Some important terms

- DNA barcoding  → Sequence variation in a single locus (e.g. ITS) in a single specimen
- Metabarcoding  → Sequence variation in a single locus (e.g. ITS) in a community
Mixed DNA templates
- Metagenomics  → Genome wide sequence variation in a community
Mixed DNA templates
- Metatranscriptomics  → cDNA sequence variation in a community
Mixed RNA

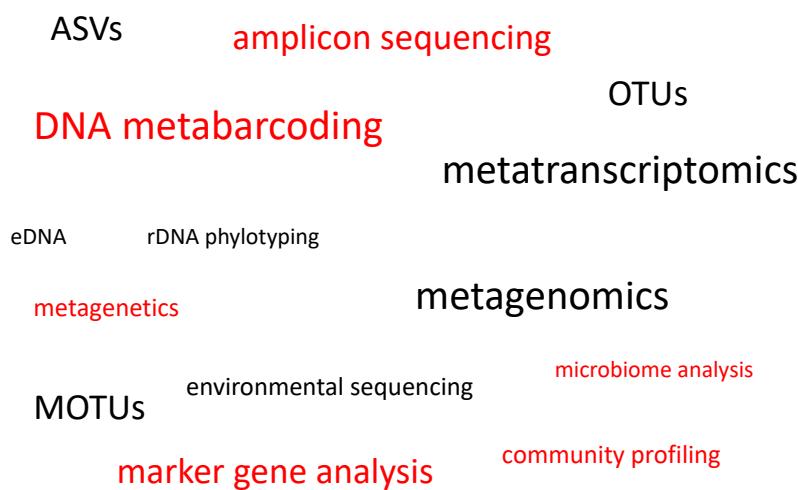
Some important terms

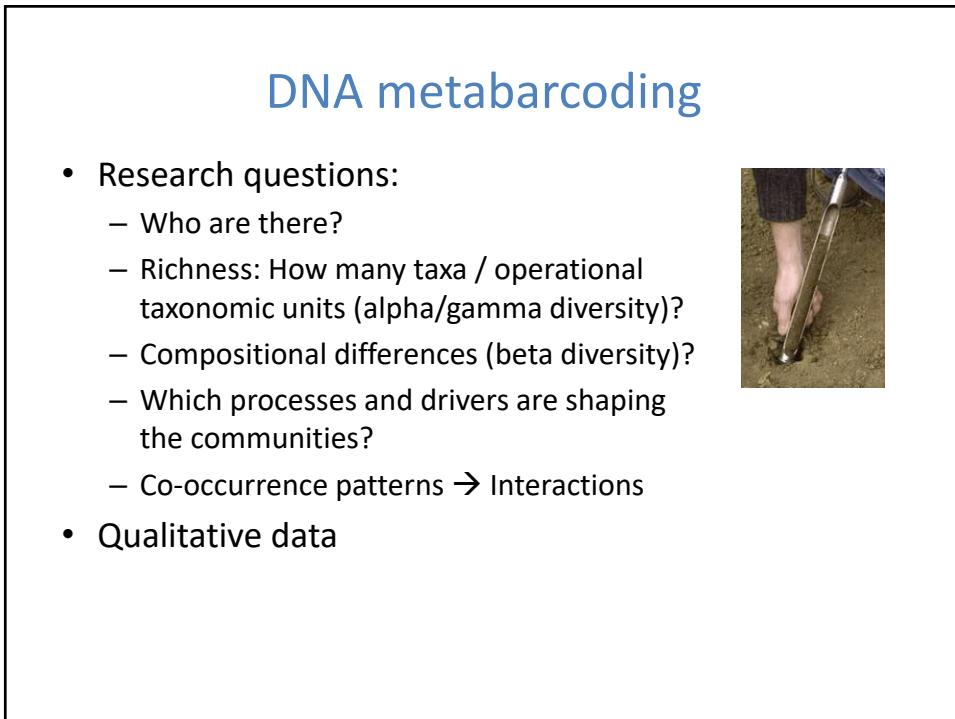
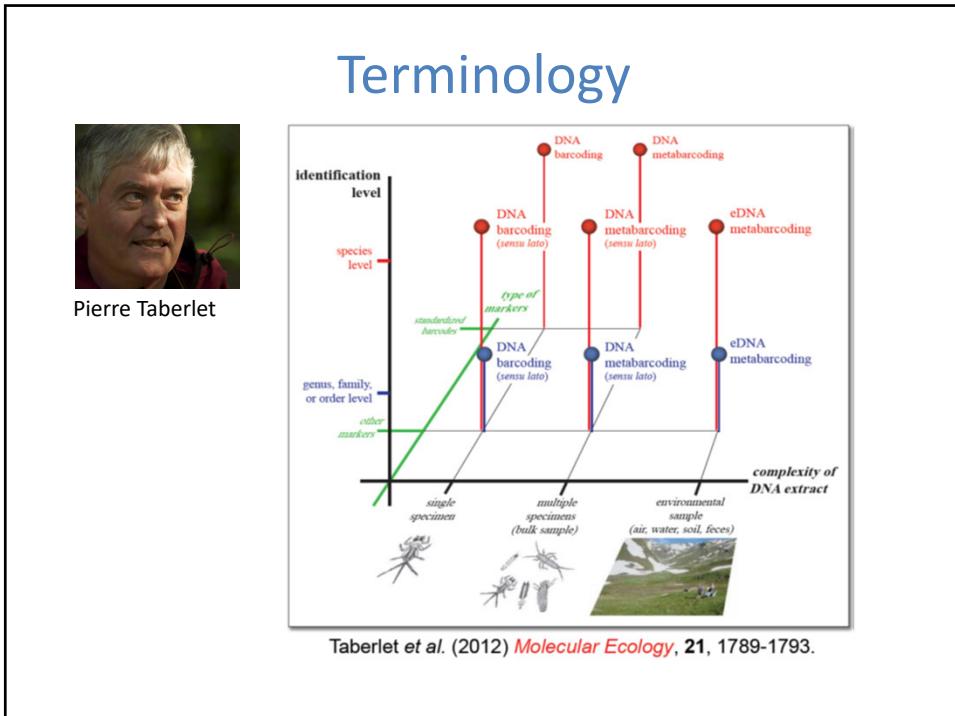
- Metabarcoding  → Sequence variation in a single locus (e.g. 16S)
Mixed DNA templates
 - Metagenomics  → Genome wide sequence variation
Mixed DNA templates
 - Metatranscriptomics  → cDNA sequence variation
Mixed RNA
- Who are there?**
- Which genes (and who) are there?**
- Who are active and doing what?**

Some important terms

- Metabarcoding  Sequence variation in a single locus (e.g. 16S) **Who are there?**
- Metagenomics  Genome wide sequence variation **Which genes (and who) are there?**
- Metatranscriptomics  cDNA sequence variation **Who are active and doing what?**

Some confusing / important terms





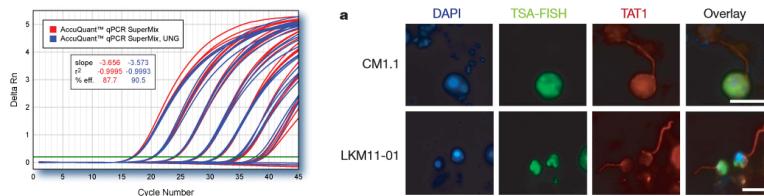
DNA metabarcoding

- Research questions:
 - Who are there?
 - Richness: How many taxa / operational taxonomic units (alpha/gamma diversity)?
 - Compositional differences (beta diversity)?
 - Which processes and drivers are shaping the communities?
 - Co-occurrence patterns → Interactions
- Qualitative data
- Quantitative data
 - Abundance: Who are common – who are rare...?



DNA metabarcoding

- Often functions as a first step looking into poorly characterized habitats and study systems
- Often generates hypotheses that can be addressed more in-depth with more taxon-specific approaches (e.g. qPCR or genomics) or genomics



DNA metabarcoding: From «wild west» towards an established approach



- Primary phase**
- Poor replication
 - Lack of controls
 - Lack of insight into important biases
 - Poor bioinformatics approaches

EDITORIAL

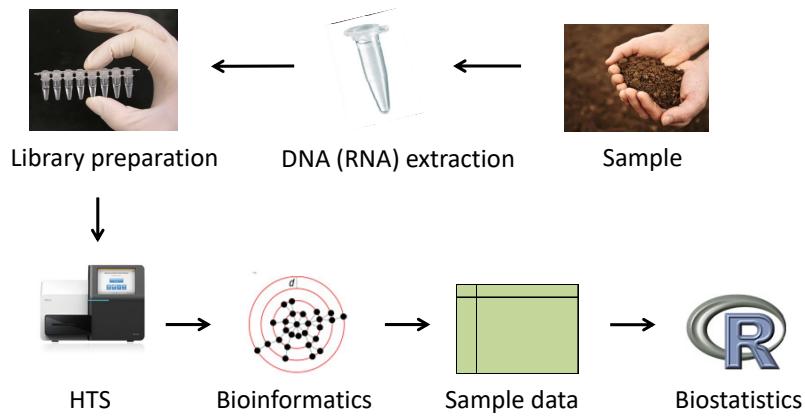
MOLECULAR ECOLOGY WILEY

DNA metabarcoding—Need for robust experimental designs to draw sound ecological conclusions
Zinger et al. 2019, Molecular Ecology Resources

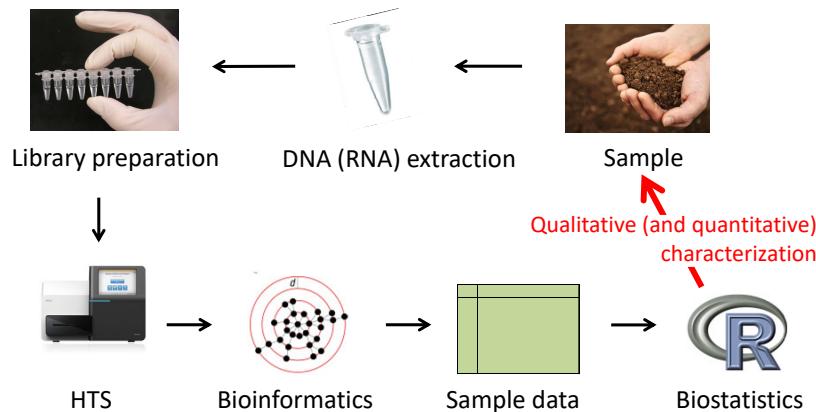


Established scientific approach with a set of widely accepted guidelines

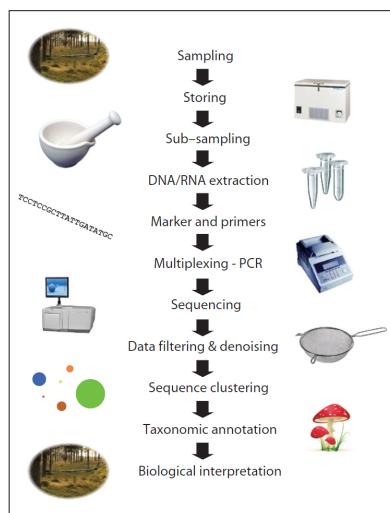
General workflow in DNA metabarcoding studies



General workflow in DNA metabarcoding studies



DNA metabarcoding - many steps



Many
steps...



... to go
wrong

Lindahl et al. 2013

EDITORIAL

MOLECULAR ECOLOGY WILEY

DNA metabarcoding—Need for robust experimental designs to draw sound ecological conclusions

Zinger et al. 2019. Molecular Ecology Resources

RESEARCH ARTICLE

Methods in Ecology and Evolution BRITISH ECOLOGICAL SOCIETY

Scrubinizing key steps for reliable metabarcoding of environmental samples

Antton Alberdi¹ | Ostaizka Aizpurua¹ | M. Thomas P. Gilbert^{1,2,3} | Kristine Bohmann^{1,4}

Alberdi et al. 2017

DNA metabarcoding - many steps

```

graph TD
    Sampling --> Storing
    Storing --> Subsampling
    Subsampling --> DNAExtraction[DNA/RNA extraction]
    DNAExtraction --> MarkerPrimers[Marker and primers]
    MarkerPrimers --> MultiplexingPCR[Multiplexing - PCR]
    MultiplexingPCR --> Sequencing
    Sequencing --> DataFiltering[Data filtering & denoising]
    DataFiltering --> SequenceClustering[Sequence clustering]
    SequenceClustering --> TaxonomicAnnotation[Taxonomic annotation]
    TaxonomicAnnotation --> BioInterpretation[Biological interpretation]
    
```

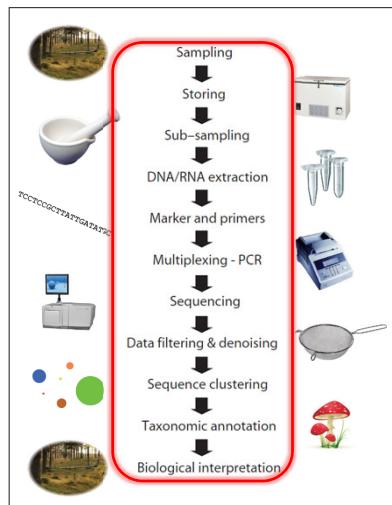
The flowchart illustrates the DNA metabarcoding process. It begins with 'Sampling' (represented by a forest icon), followed by 'Storing' (represented by a beaker icon). This leads to 'Sub-sampling' (represented by a small beaker icon) and then 'DNA/RNA extraction' (represented by a test tube icon). Next is 'Marker and primers' (represented by a PCR machine icon), then 'Multiplexing - PCR' (represented by a computer monitor icon). The process continues with 'Sequencing' (represented by a sequencing instrument icon), then 'Data filtering & denoising' (represented by a pan icon), 'Sequence clustering' (represented by a cluster icon), 'Taxonomic annotation' (represented by a mushroom icon), and finally 'Biological interpretation' (represented by a forest icon).

Many steps...

... to go wrong

Lindahl et al. 2013

DNA metabarcoding - many steps



Lindahl et al. 2013

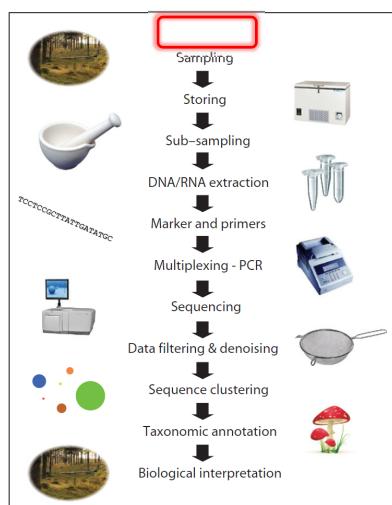


Many steps...



... to go wrong

DNA metabarcoding - many steps



Lindahl et al. 2013



Many steps...



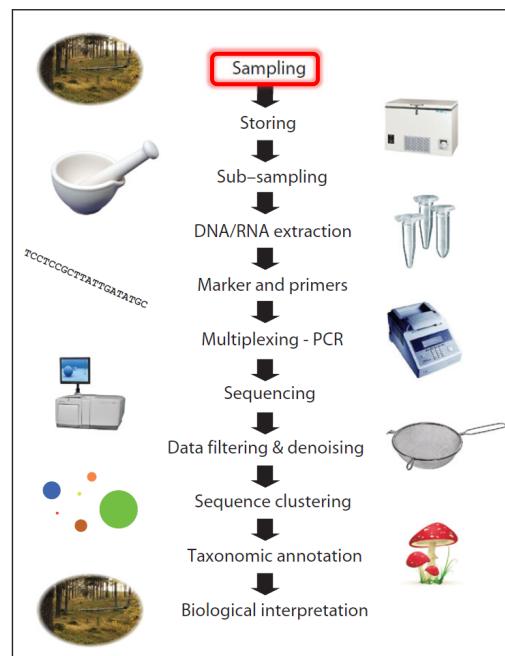
... to go wrong

If new study system – conduct a pilot?

- Which sampling scheme?
- How many replicates?
- Which extraction protocol
- Which primers?
- Which sequencing depth?
- Which sequencing technique?
- Etc.



→ Depends on the alpha, beta and gamma diversity, that you might not know anything about..



Representativeness (in space)

- Many communities highly heterogenous
- Should obtain samples that are representative
- They should not be too small
- If you are not interested in the small scale variation in itself → pool sub-samples?



Fig. 5. What is the optimal relationship between primary sample size and the analytical sample volume (insert) and how can it come about? When sample size increases one can intuitively understand that the sample becomes more representative. But at the same time, today's analytical volumes continue to decrease (insert) as the analytical instruments become more and more precise. For all heterogeneous materials, there is consequently an intrinsic contradiction between primary sampling representativity and the instrumental analytical volume requirements. This is the root cause of all sampling and representativity issues.

Representativeness (in time)

- Many communities often display high temporal variation! Repeated temporal sampling?



Examples: Insects
and fungal spores
in the air

«Replicate or lie»

Environmental Microbiology (2010) 12(7), 1806–1810
doi:10.1111/j.1462-2920.2010.02194.x

Opinion
Replicate or lie

James P. Prosser¹
Institute of Biological and Environmental Sciences,
University of Aberdeen, Crookston Building, St.
Mary's Street, Aberdeen, AB22 2PU, UK

Introduction
Anderson and colleagues (2005) recently published a paper in this journal that has received considerable attention and interest during the early years of research, but has also been widely cited in the media. The paper concluded that highly likely most of the diversity of microorganisms have been missed by sequencing a single sample. This was based primarily on statistical analysis and its significance was overestimated. I am not going to discuss the statistics in this paper, and my authority on the consequences, and the lack of them, of this paper is limited to the issue of microbial diversity. The authors are positive with species richness and diversity, and negative with species evenness and evenness of distribution. The authors conclude that the new sequencing technologies (e.g. metagenomics) are high-throughput sequencing methods that can overcome the limitations of previous methods beyond these techniques and boosted studies of microbial diversity.

Why replicate?
The main point of the journal could be that the article describes and discusses applications of statistical analysis in environmental microbiology. The main point of the paper is that the base and fundamental aspect – the need for replication – is ignored. In this paper, the authors did not even consider writing to compare bacterial abundance in two samples. They did not even consider the effect of sequencing depth. They did not even consider the effect of sequencing error rate. Only the effect of sequencing error rate is greater than one like for the other. They did not even consider the effect of the number of samples. Their basic mistake includes lack of and/or variability among them.

Received 7 December 2009; accepted 2 January 2010. *Correspondence to: James P. Prosser, Institute of Biological and Environmental Sciences, University of Aberdeen, St Mary's Street, Aberdeen, AB22 2PU, UK.
© 2010 Society for Applied Microbiology and Blackwell Publishing Ltd

Clone library analysis and pyrosequencing

	Number of articles	% with replicates
<i>Appl Environ Microbiol</i>	60	23
<i>Environ Microbiol</i>	47	15
<i>FEMS Microbiol Ecol</i>	29	24
<i>ISME J</i>	23	13
<i>Microbial Ecol</i>	22	9
Total	181	18

It doesn't help that you
are dealing with HTS data if
you don't replicate properly!

Prosser JI. 2010, *Environmental Microbiology*

Received 4 October 2017 | Revised 10 May 2018 | Accepted 14 May 2018
DOI: 10.1111/1365-2745.13297

INVITED TECHNICAL REVIEW **WILEY** **MOLECULAR ECOLOGY**

Towards robust and repeatable sampling methods in eDNA-based studies

Ian A. Dickie^{1,2} | Stéphane Boyer^{3,4} | Hannah L. Buckley⁵ | Richard P. Duncan⁶ | Paul P. Gardner⁷ | Ian D. Hogg^{7,8} | Robert J. Holdaway⁹ | Gavin Lear¹⁰ | Andreas Makiola¹ | Sergio E. Morales¹¹ | Jeff R. Powell¹² | Louise Weaver¹³

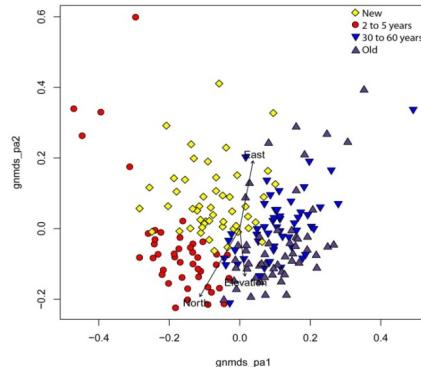
Abstract

DNA-based techniques are increasingly used for measuring the biodiversity (species presence, identity, abundance and community composition) of terrestrial and aquatic ecosystems. While there are numerous reviews of molecular methods and bioinformatic steps, there has been little consideration of the methods used to collect samples upon which these later steps are based. This represents a critical knowledge gap, as methodologically sound field sampling is the foundation for subsequent analyses. We reviewed field sampling methods used for metabarcoding studies of both terrestrial and freshwater ecosystem biodiversity over a nearly three-year period ($n = 75$). We found that 95% ($n = 71$) of these studies used subjective sampling methods and inappropriate field methods and/or failed to provide critical methodological information. It would be possible for researchers to replicate only 5% of the metabarcoding studies in our sample, a poorer level of reproducibility than for ecological studies in general. Our findings suggest greater attention to field sampling methods, and reporting is necessary in eDNA-based studies of biodiversity to ensure robust outcomes and future reproducibility. Methods must be fully and accurately reported, and protocols developed that minimize subjectivity. Standardization of sampling protocols would be one way to help to improve reproducibility and have additional benefits in allowing compilation and comparison of data from across studies.

Biological replicates



Biological replicates are biologically distinct samples which show biological variation



Fungal communities associated with mosses in different forest management types

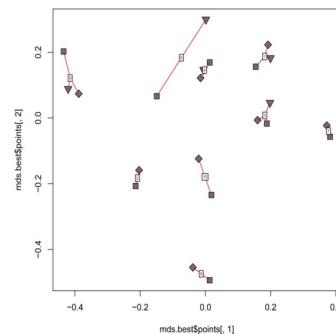
Davey et al. 2014. FEMS Microbial Ecology

Technical replicates

- Some samples should/can be analyzed multiple times
- Reveals the variability (experimental error) of the analysis → allows to set limits for what is meaningful and significant data
- How important?

Communities with low DNA content!!

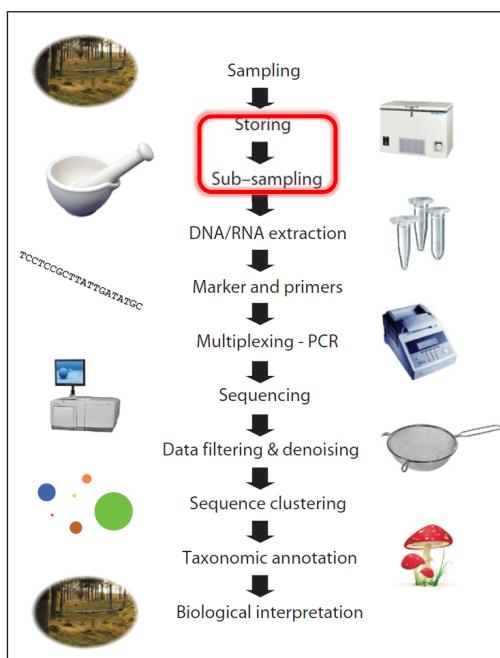
Technical replicates are repeated measurements of a sample, which show variation of the measuring equipment and protocol



Davey et al. 2014. FEMS Microbial Ecology

Sample types

1. Biological replicates
2. Technical replicates



Storage

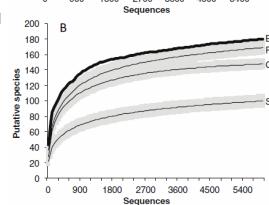
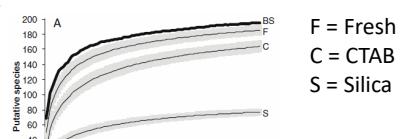
- Unappropriate storage may introduce severe biases!
- Community members can respond quickly to altered conditions
- 'Arrest' the communities!
- Process the samples asap. If needed, long time storage at -80C often suggested



Storage



U'Ren et al. 2014. Tissue storage and primer selection influence pyrosequencing-based inferences of diversity and community composition of endolichenic and endophytic fungi. Mol Ecol Res.



Sub-sampling and homogenization

- Protocols for nucleic acid extraction are normally based on small amounts
- Field samples are often much larger, and careful homogenization of material is required to reduce the sample to a smaller but still representative subsamples
- The most commonly used techniques are bead beating and/or crushing in liquid nitrogen.



Homogenization of samples

- After homogenization, new spatial structures may easily be created in the samples, for example by density fractionation at the slightest bumping.
- To make proper comparisons: DNA should be extracted from equivalent amounts of starting material!

MOLECULAR ECOLOGY

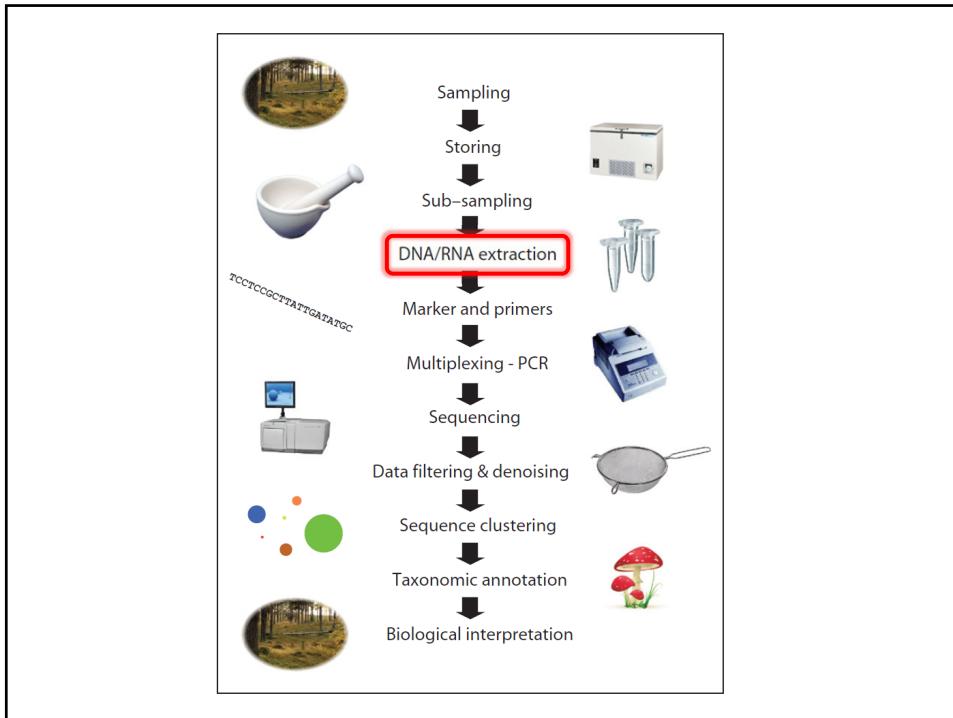
Molecular Ecology (2012) 21, 1816–1820

doi: 10.1111/j.1365-294X.2011.05317.x

**Soil sampling and isolation of extracellular DNA
from large amount of starting material suitable
for metabarcoding studies**

PIERRE TABERLET, SOPHIE M. PRUD'HOMME, ETIENNE CAMPIONE, JULIEN ROY,
CHRISTIAN MIQUEL, WASIM SHEHZAD, LUDOVIC GIELLY, DELPHINE RIOUX,
PHILIPPE CHOLER, JEAN-CHRISTOPHE CLÉMENT, CHRISTELLE MELODELIMA,
FRANÇOIS POMPANON and ERIC COISSAC

Laboratoire d'Ecologie Alpine, CNRS UMR 5553, Université Joseph Fourier, BP 53, F-38041 Grenoble Cedex 9, France



DNA extraction

An evaluation of commercial DNA extraction kits for the isolation of bacterial spore DNA from soil
S.M. Dineen^{1,2}, R. Aranda IV^{1,2}, D.L. Anders³ and J.M. Robertson²

¹ Visiting Scientist, Federal Bureau of Investigation Laboratory, Quantico, VA, USA
² Counterterrorism and Foreign Science Research Unit, Federal Bureau of Investigation Laboratory, Quantico, VA, USA
³ Hazardous Materials Science Response Unit, Federal Bureau of Investigation Laboratory, Quantico, VA, USA

Influence of DNA extraction and PCR amplification on studies of soil fungal communities based on amplicon sequencing
Lihui Xu, Sabine Ravanska, John Larsson, and Magnus Nilssen¹

Molecular biology, genetics and biotechnology
Effect of DNA extraction and sample preservation method on rumen bacterial population
Katerina Fliegerova^{a,b,*}, Ilma Tapiö^b, Aurelie Bonin^c, Jakub Mrazeck^a, Maria Luisa Callegari^c, Paolo Ratti^c, Alireza Bayat^d, Johanna Vilki^c, Jan Kopečný^a, Kevin J. Shilling^{a,b}, Frederic Boyer^c, Eric Coissac^c, Pierre Taberlet^c, R. John Wallace^c, Pierre Chadoeuf^c, and Sébastien Lepage^a

Effect of DNA Extraction Methods and Sampling Techniques on the Apparent Structure of Cow and Sheep Rumen Microbial Communities
Gemma Henderson¹, Faith Cox¹, Sanna Kuitiainen¹, Vahideh Heidarian Miri², Michael Zehfouj², Samantha J. Noel³, Garry C. Waghorn², Peter H. Jansson¹

The Impact of Different DNA Extraction Kits and Laboratories upon the Assessment of Human Gut Microbiota Composition by 16S rRNA Gene Sequencing
Nicholas A. Kennedy¹, Alan W. Walker², Susan H. Berry³, Sylvia H. Duncan⁴, Freda M. Farquharson⁴, Petra Louis⁵, John M. Thomson⁵, UK IBD Genetics Consortium⁶, Jack Satsangi¹, Harry J. Flint⁶, Julian Parkhill², Charlie W. Lee^{1*}, Georgina L. Hold²

* indicates author for correspondence

DNA extraction

- Should yield high and uniform amounts of DNA
- Concentration of PCR inhibitors minimized
- Same protocol for all samples!
- If no proper literature are available on your study system → conduct a pilot?!

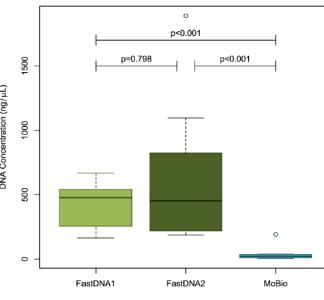
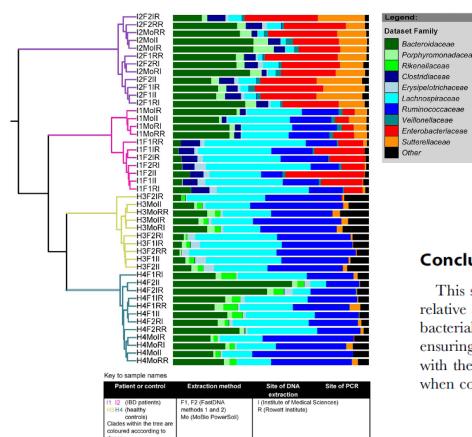


- MoBio Power Soil?
- FastDNA kit for Soil?
- EZNA Soil kit?
- CTAB + cleanup kit?

DNA extraction

The Impact of Different DNA Extraction Kits and Laboratories upon the Assessment of Human Gut Microbiota Composition by 16S rRNA Gene Sequencing

Nicholas A. Kennedy¹, Alan W. Walker², Susan H. Berry², Sylvia H. Duncan³, Freida M. Farquharson⁴, Petra Louis⁵, John M. Thomson⁵, UK IBD Genetic Consortium⁶, Jack Satsangi¹, Harry J. Flint⁴, Julian Parkhill⁷, Charlie W. Lee³, Georgina L. Hold^{1,*}



Conclusions

This study demonstrates important differences in the yield and relative abundance of key bacterial families for kits used to isolate bacterial DNA from stool. This highlights the importance of ensuring that all samples to be analyzed together are prepared with the same DNA extraction method, and the need for caution when comparing studies that have used different methods.

Note: Lack of true replicates

DNA extraction

- Both extracellular and intracellular DNA are normally co-extracted.
- Method for extraction of extracellular DNA from large amount of starting material:

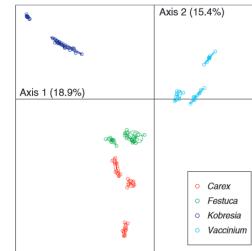
MOLECULAR ECOLOGY

Molecular Ecology (2012) 21, 1816–1820

doi: 10.1111/j.1365-294X.201

Soil sampling and isolation of extracellular DNA from large amount of starting material suitable for metabarcoding studies

PIERRE TABERLET, SOPHIE M. PRUD'HOMME, ETIENNE CAMPIONE, JULIEN ROY, CHRISTIAN MIQUEL, WASIM SHEHZAD, LUDOVIC GIELLY, DELPHINE RIOUX, PHILIPPE CHOLER, JEAN-CHRISTOPHE CLÉMENT, CHRISTELLE MELODELIMA, FRANÇOIS POMPANON and ERIC COISSAC
Laboratoire d'Ecologie Alpine, CNRS UMR 5553, Université Joseph Fourier, BP 53, F-38041 Grenoble Cedex 9, France



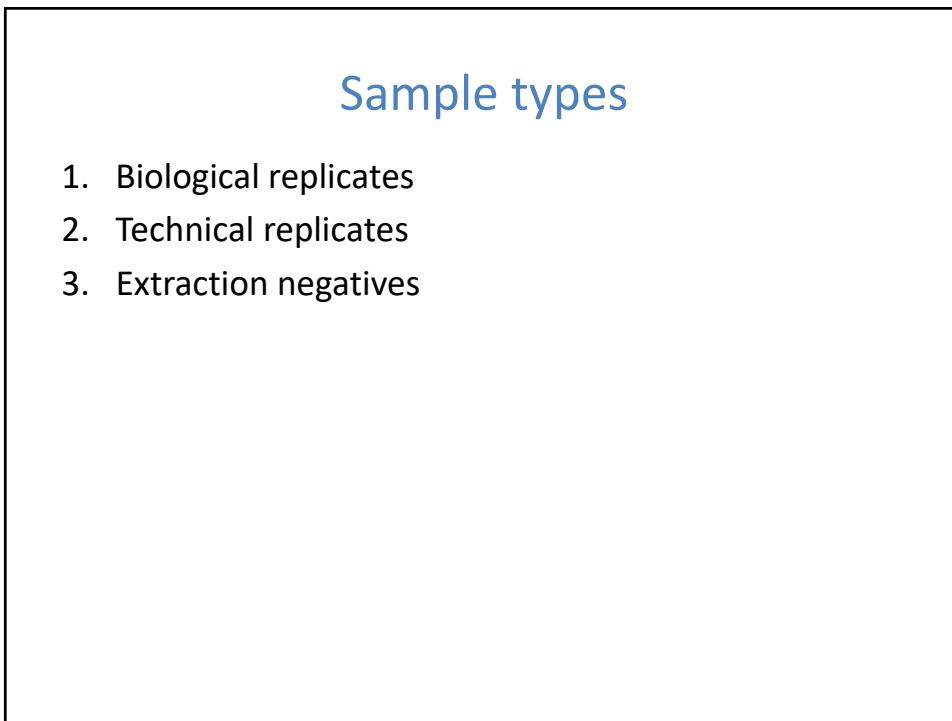
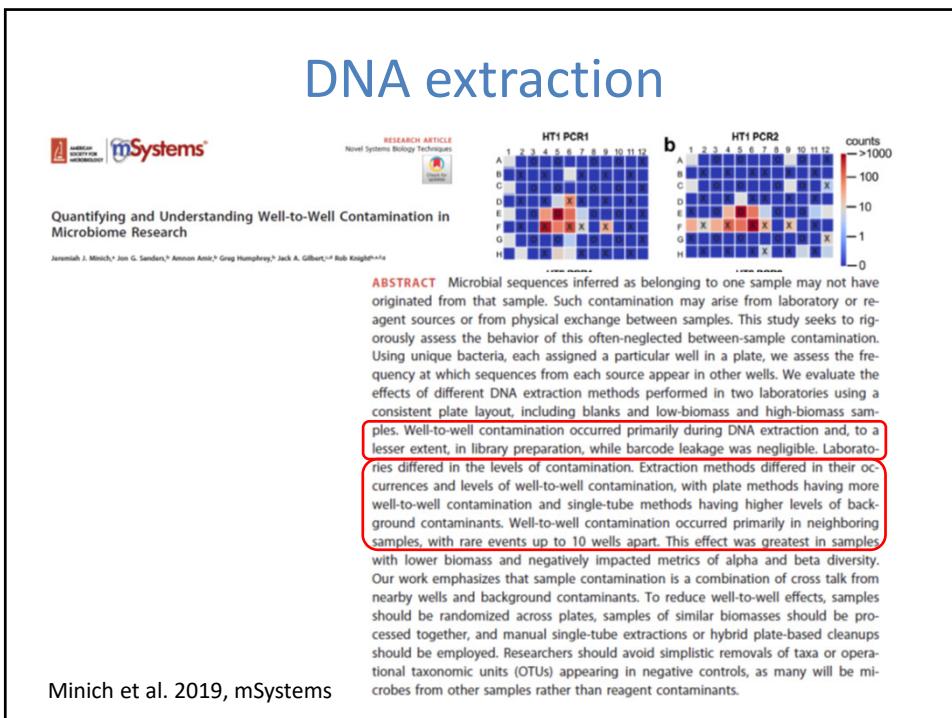
Replicates!

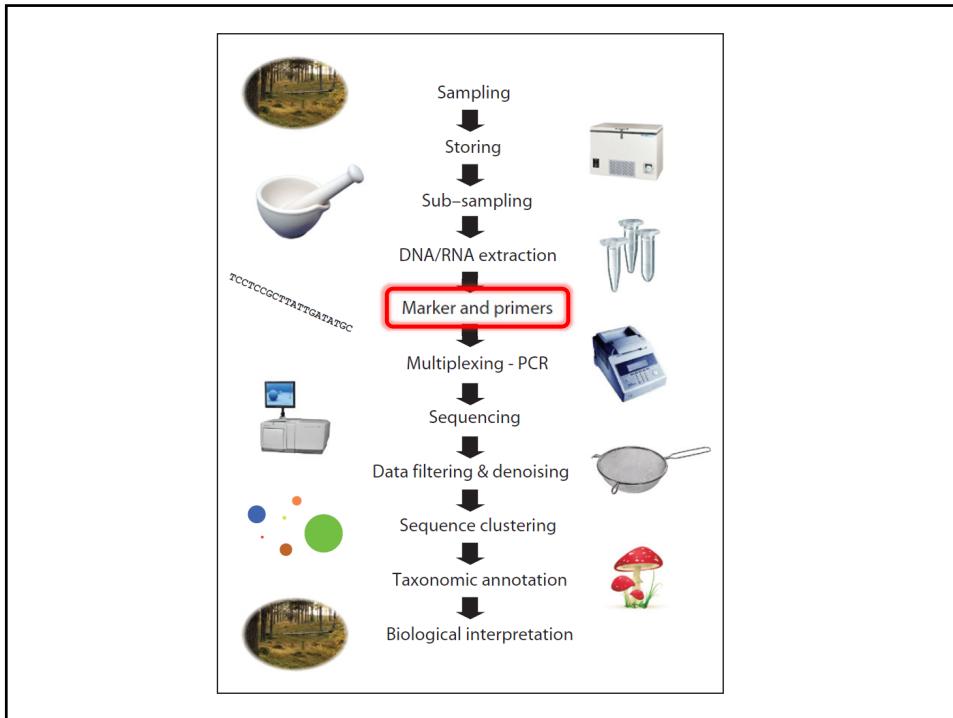
DNA extraction

- Should yield high and uniform amounts of DNA
- Concentration of PCR inhibitors minimized
- Same protocol for all samples!
- If no proper literature are available on your study system → conduct a pilot?!
- Extraction negatives!



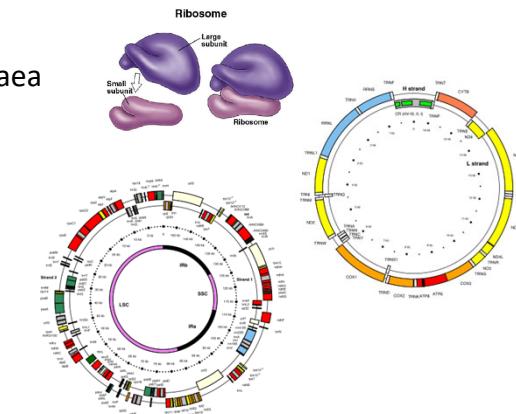
- MoBio Power Soil?
- FastDNA kit for Soil?
- EZNA Soil kit?
- CTAB + cleanup kit?





Markers used in DNA metabarcoding

- Standard markers (<500 bp):
 - 18S: Eukaryotes
 - 16S: Bacteria/archaea
 - ITS: Fungi & plants
 - COI: Metazoa
 - *RbcL*: Plants
 - *trnL*: Plants

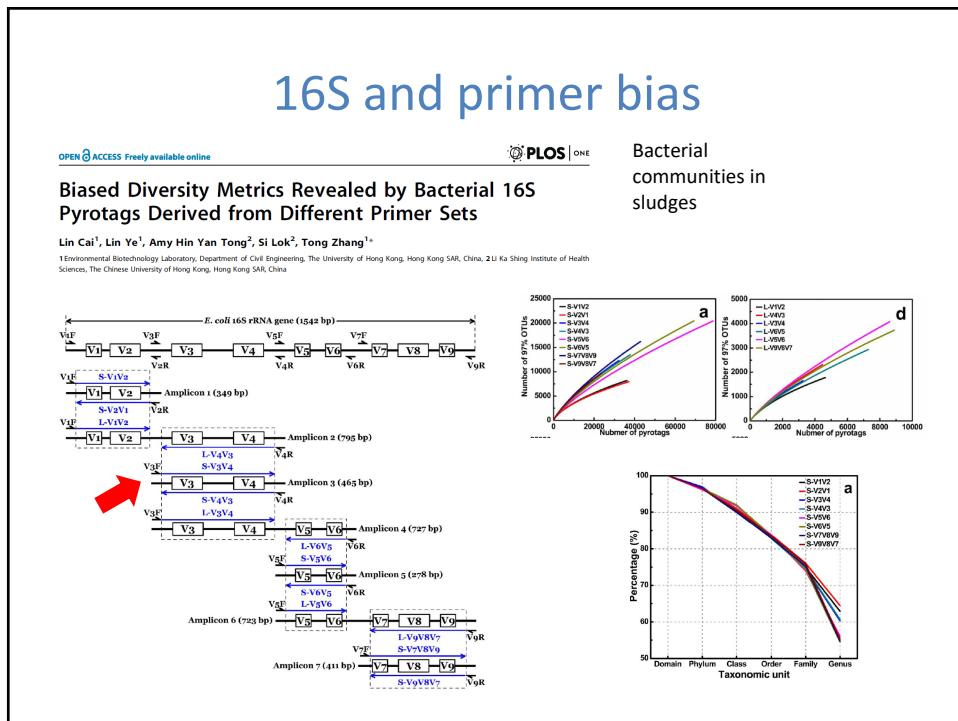
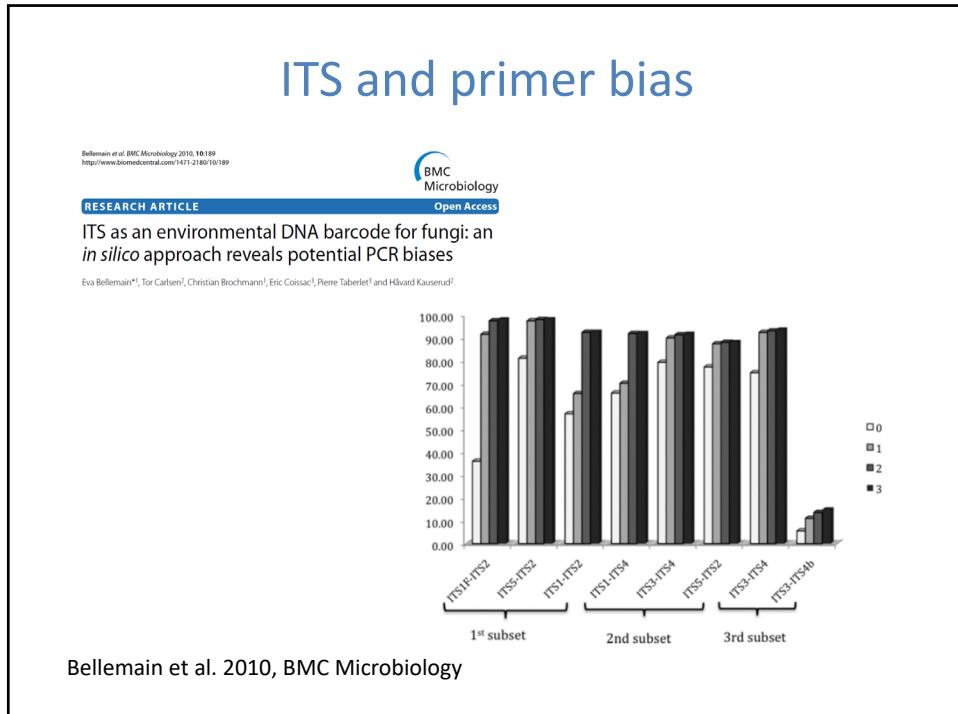


Markers in DNA metabarcoding

- The ideal marker should:
 - Have primer sites that are shared by all target organisms
 - Be easy to amplify
 - Be of appropriate length for efficient amplification and sequencing
 - Be of similar length
 - No intragenomic variation (i.e. no paralogs)
 - Similar number of copies
 - Be possible to align
 - Have high interspecific variation
 - Have low intraspecific variation
- No known markers meet all these requirements!

Markers in DNA metabarcoding

- The ideal marker should:
 - Have primer sites that are shared by all target organisms
 - Be easy to amplify
 - Be of appropriate length for efficient amplification and sequencing
 - Be of similar length
 - No intragenomic variation (i.e. no paralogs)
 - Similar number of copies
 - Be possible to align
 - Have high interspecific variation
 - Have low intraspecific variation
- No known markers meet all these requirements!



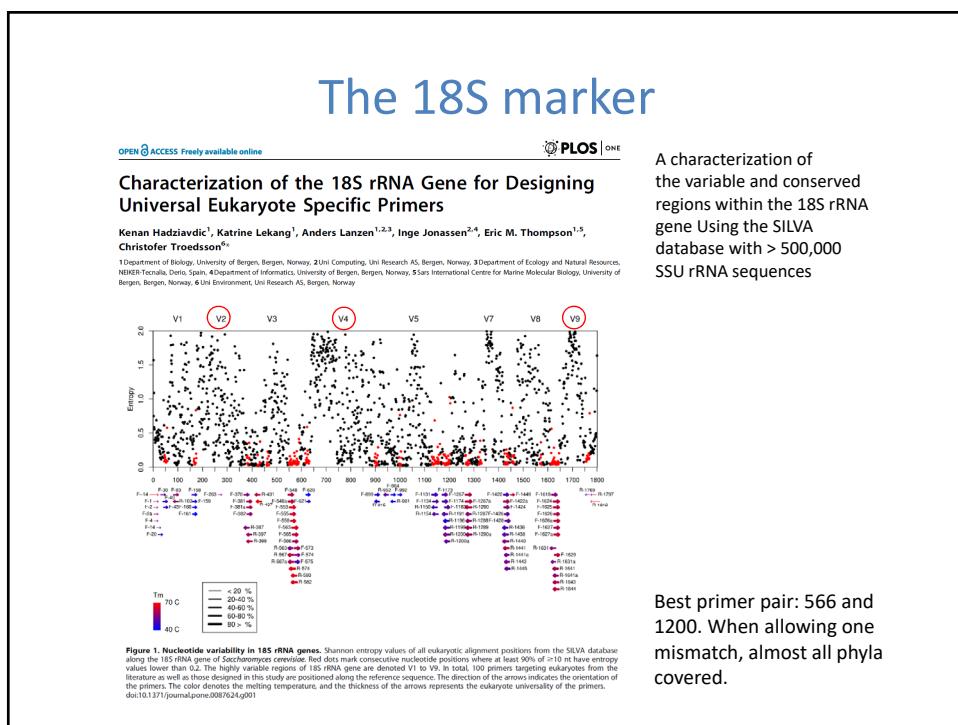
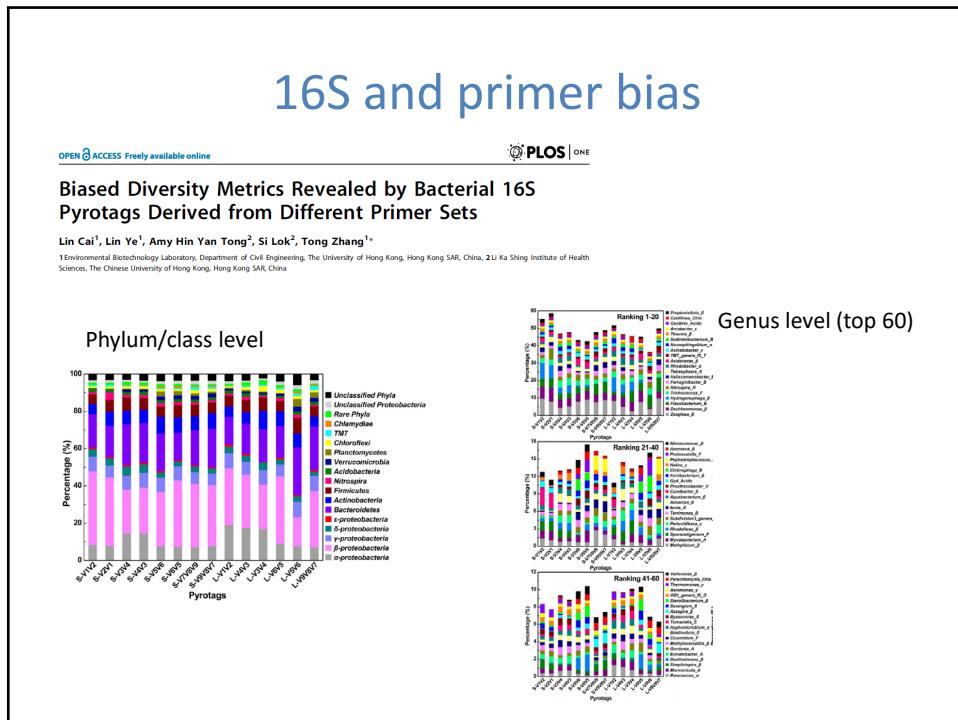


Figure 1. Nucleotide variability in 18S rRNA genes. Shannon entropy values of all eukaryotic alignment positions from the SILVA database along the 18S rRNA gene of *Saccharomyces cerevisiae*. Red dots mark consecutive nucleotide positions where at least 90% of >10 nt have entropy values lower than 0.2. The highly variable regions of 18S rRNA gene are denoted V1 to V9. In total, 100 random target eukaryotes from the SILVA database were used in this study. The arrows indicate the annealing positions of the primers. The length of the arrow indicates the coverage of the primers. The color denotes the melting temperature, and the thickness of the arrows represents the eukaryote universality of the primers.

Markers in DNA metabarcoding

- The ideal marker should:
 - Have primer sites that are shared by all target organisms
 - Be easy to amplify
 - Be of appropriate length for efficient amplification and sequencing
 - Be of similar length
 - No intragenomic variation (i.e. no paralogs)
 - Be possible to align
 - Have high interspecific variation
 - Have low intraspecific variation
- No known markers meet all these requirements!

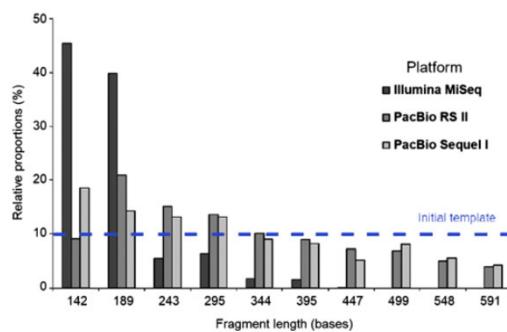
With length variation
 → length biases
 introduced during both
 PCR and sequencing!



Methods

Optimized metabarcoding with Pacific biosciences enables semi-quantitative analysis of fungal communities

Carles Castaño¹, Anna Berlin¹, Mikael Brandström Durling¹, Katharina Ihrmark¹, Björn D. Lindahl², Jan Stenlid¹, Karina E. Clemmensen^{1*} and Ake Olson^{1*}



Markers in DNA metabarcoding

- The ideal marker should:
 - Have primer sites that are shared by all target organisms
 - Be easy to amplify
 - Be of appropriate length for efficient amplification and sequencing
 - Be of similar length
 - **No intragenomic variation (i.e. no paralogs)**
 - Similar number of copies
 - Be possible to align
 - Have high interspecific variation
 - Have low intraspecific variation
- No known markers meet all these requirements!

(Intra)genomic variability in 16S

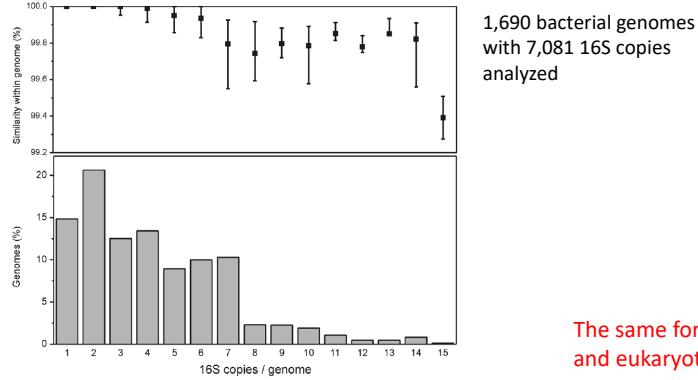
OPEN  ACCESS Freely available online

 PLOS ONE

The Variability of the 16S rRNA Gene in Bacterial Genomes and Its Consequences for Bacterial Community Analyses

Tomáš Větrovský, Petr Baldrian*

Laboratory of Environmental Microbiology, Institute of Microbiology of the Academy of Sciences of the Czech Republic, Praha, Czech Republic



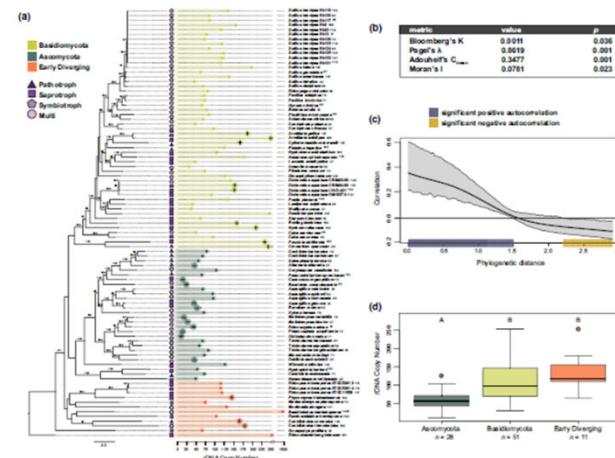
Markers in DNA metabarcoding

- The ideal marker should:
 - Have primer sites that are shared by all target organisms
 - Be easy to amplify
 - Be of appropriate length for efficient amplification and sequencing
 - Be of similar length
 - No intragenomic variation (i.e. no paralogs)
 - **Similar number of copies**
 - Be possible to align
 - Have high interspecific variation
 - Have low intraspecific variation
- No known markers meet all these requirements!

ORIGINAL ARTICLE WILEY MOLECULAR ECOLOGY

Genome-based estimates of fungal rDNA copy number variation across phylogenetic scales and ecological lifestyles

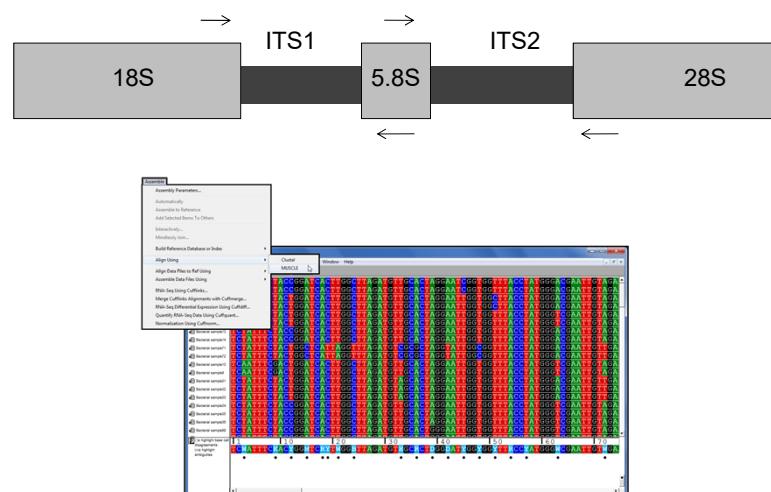
Lotus A. Lofgren¹ | Jessie K. Uehling² | Sara Branco³ | Thomas D. Bruns² | Francis Martin⁴ | Peter G. Kennedy^{1,5}



Markers in DNA metabarcoding

- The ideal marker should:
 - Have primer sites that are shared by all target organisms
 - Be easy to amplify
 - Be of appropriate length for efficient amplification and sequencing
 - Be of similar length
 - No intragenomic variation (i.e. no paralogs)
 - Similar number of copies
 - **Be possible to align (in a multiple alignment)**
 - Have high interspecific variation
 - Have low intraspecific variation
- No known markers meet all these requirements!

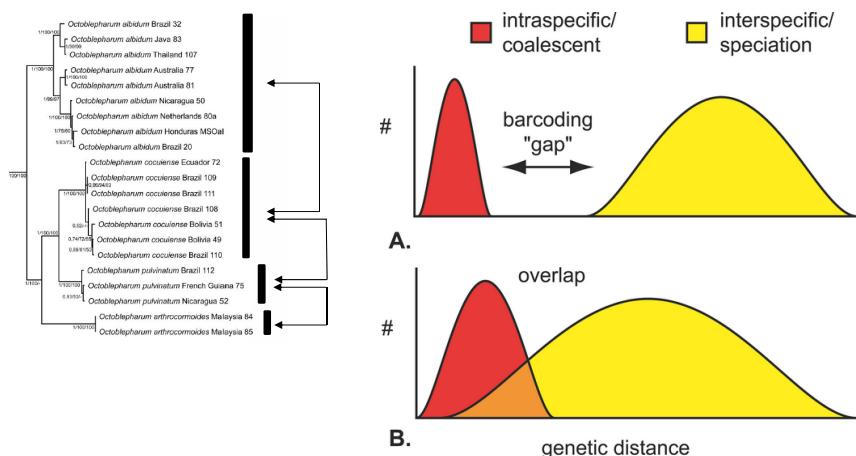
Possible to align



Markers in DNA metabarcoding

- The ideal marker should:
 - Have primer sites that are shared by all target organisms
 - Be easy to amplify
 - Be of appropriate length for efficient amplification and sequencing
 - Be of similar length
 - No intragenomic variation (i.e. no paralogs)
 - Be possible to align
 - Possess high interspecific variation
 - Possess low intraspecific variation
- No known markers meet all these requirements!

The barcoding gap



How conserved/variable are the marker?

- 18S and 16S: Low variability, low intraspecific variation, low interspecific variation
- ITS: High variability, high intraspecific variation, high 'interspecific' variation



- Impact how the bioinformatics analyses should be conducted (i.e. no single way)!!

Multiple markers/primers?

Drummond et al. *GigaScience* (2015) 4:46
DOI 10.1186/s13742-015-0086-1



RESEARCH

Open Access



Evaluating a multigene environmental DNA approach for biodiversity assessment

Alexei J. Drummond^{1,2*}, Richard D. Newcomb^{1,3,4}, Thomas R. Buckley^{1,3,5}, Dong Xie^{1,2}, Andrew Dopheide^{1,3,4}, Benjamin CM Potter^{1,3}, Joseph Heled^{1,2}, Howard A. Ross^{1,3}, Leah Tooman^{1,4}, Stefanie Grosser^{1,5}, Duckchul Park⁵, Nicholas J. Demetras⁶, Mark I. Stevens^{6,7}, James C. Russell^{1,3,9}, Sandra H. Anderson³, Anna Carter^{1,10} and Nicola Nelson^{1,10}

- Inherent problem in complex samples: Difficult to link alleles/variants across DNA markers/loci
- Often lack proper reference databases for multiple markers.

Long-read metabarcoding

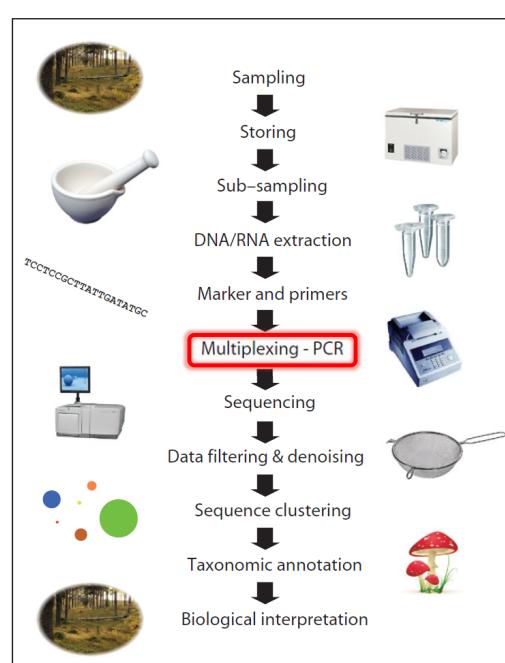
RESOURCE ARTICLE

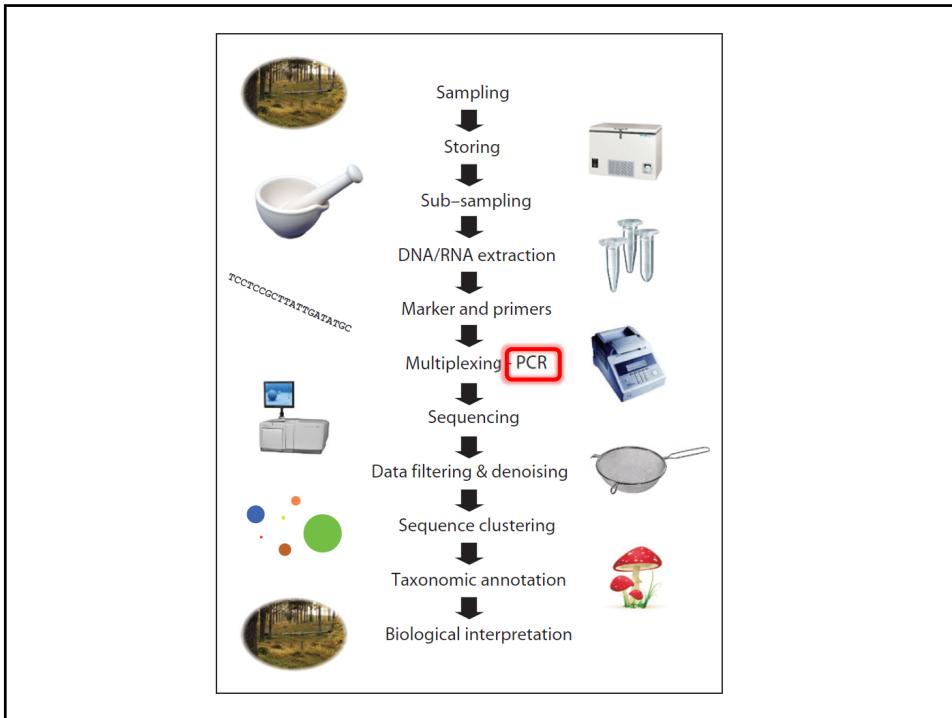
MOLECULAR ECOLOGY
RESOURCES

Long-read metabarcoding of the eukaryotic rDNA operon to phylogenetically and taxonomically resolve environmental diversity

Mahwash Jamy¹  | Rachel Foster² | Pierre Barbera³ | Lucas Czech³ | Alexey Kozlov³ | Alexandros Stamatakis^{3,4} | Gary Bending⁵ | Sally Hilton⁵ | David Bass^{2,6}  | Fabien Burki¹

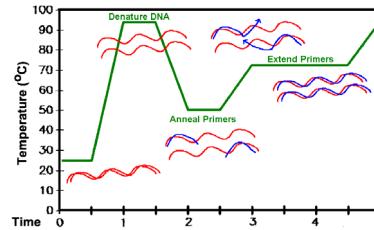
- Phylogenetic framework → more secure taxonomic placement
- Comes with extra challenges
 - Harder to amplify
 - More chimeric sequences
 - Lower sequencing depth (PacBio or Nanopore)





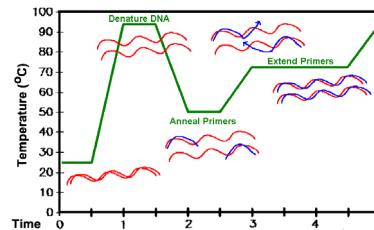
PCR

- Different relevant factors during PCR:
 - Which polymerase enzyme (proofreading or not)?
 - Which RAMP speed?
 - How many cycles?
 - Which annealing temperature?
 - Multiple/replicate PCR reactions?
 - PCR negatives!



PCR

- Different relevant factors during PCR:
 - Which polymerase enzyme (proofreading or not)?
 - Which RAMP speed?
 - How many cycles?
 - Which annealing temperature?
 - Multiple/replicate PCR reactions?
 - PCR negatives!



Polymerase enzyme

The Journal of Microbiology (2012) Vol. 50, No. 6, pp. 1071–1074
Copyright © 2012, The Microbiological Society of Korea

DOI 10.1007/s12275-012-3642-z

NOTE

Effects of PCR Cycle Number and DNA Polymerase Type on the 16S rRNA Gene Pyrosequencing Analysis of Bacterial Communities^a

Jae Hyung Ahn, Byung-Yong Kim,
Jackyeong Song, and Hang-Yeon Weon*

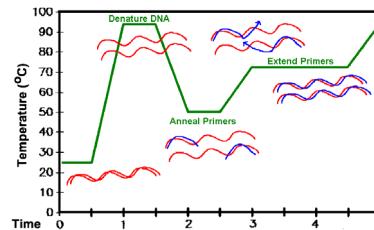
^aAgricultural Microbiology Division, National Academy of Agricultural Science, Rural Development Administration, Suwon 441-707, Republic of Korea

(Received November 19, 2012 / Accepted December 11, 2012)

The bacterial richness was overestimated at increased PCR cycle number mostly due to the occurrence of chimeric sequences, and this was more serious with a DNA polymerase having proofreading activity than with *Taq* DNA polymerase. These results suggest that PCR cycle number must be kept as low as possible for accurate estimation of bacterial richness and that particular care must be taken when a DNA polymerase having proofreading activity is used.

PCR

- Different relevant factors during PCR:
 - Which polymerase enzyme (proofreading or not)?
 - Which RAMP speed?
 - **How many cycles?**
 - Which annealing temperature?
 - Multiple/replicate PCR reactions?
 - PCR negatives!



Polymerase enzyme

The Journal of Microbiology (2012) Vol. 50, No. 6, pp. 1071–1074
Copyright © 2012, The Microbiological Society of Korea

DOI 10.1007/s12275-012-3642-z

NOTE

Effects of PCR Cycle Number and DNA Polymerase Type on the 16S rRNA Gene Pyrosequencing Analysis of Bacterial Communities^a

Jae Hyung Ahn, Byung-Yong Kim,
Jackyeong Song, and Hang-Yeon Weon*

^aAgricultural Microbiology Division, National Academy of Agricultural Science, Rural Development Administration, Suwon 441-707, Republic of Korea

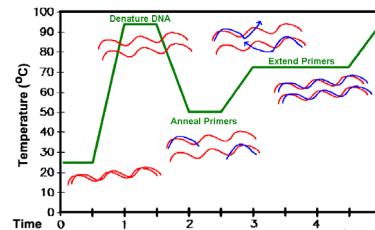
(Received November 19, 2012 / Accepted December 11, 2012)

The bacterial richness was overestimated at increased PCR cycle number mostly due to the occurrence of chimeric sequences, and this was more serious with a DNA polymerase having proofreading activity than with *Taq* DNA polymerase. These results suggest that PCR cycle number must be kept as low as possible for accurate estimation of bacterial richness and that particular care must be taken when a DNA polymerase having proofreading activity is used.

Keep the number of cycles
as low as possible!

PCR

- Different relevant factors during PCR:
 - Which polymerase enzyme (proofreading or not)?
 - Which RAMP speed?
 - How many cycles?
 - Which annealing temperature?
 - **Multiple/replicate PCR reactions?**
 - PCR negatives!



PCR replicates?

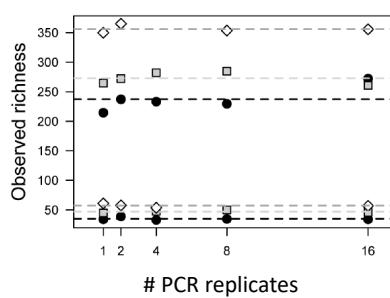
OPEN ACCESS Freely available online

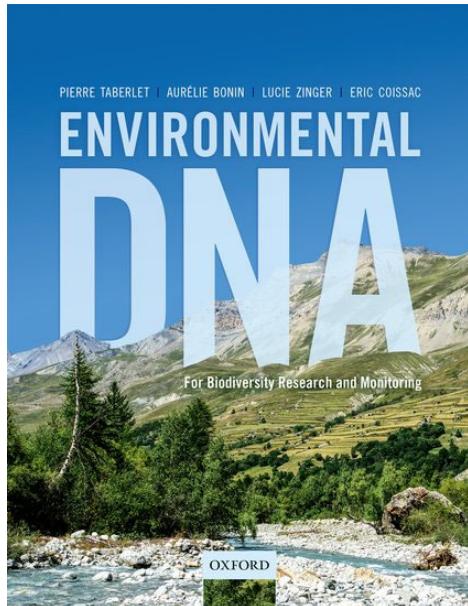
PLOS ONE

Sequence Depth, Not PCR Replication, Improves Ecological Inference from Next Generation DNA Sequencing

Dylan P. Smith, Kabir G. Peay*

Department of Biology, Stanford University, Stanford, California, United States of America





Taberlet et al: Multiple PCR amplifications and remove outliers.

Probably very context dependent: Low DNA concentrations means that stochasticity during PCR is higher → The need for multiple amplicons is higher

Sample types

1. Biological replicates
2. Technical replicates
3. Extraction negatives

+ PCR replicates: A type of technical replicates

PCR

- Different relevant factors during PCR:
 - Which polymerase enzyme (proofreading or not)?
 - Which RAMP speed?
 - How many cycles?
 - Which annealing temperature?
 - Multiple/replicate PCR reactions?
 - **PCR negatives!**

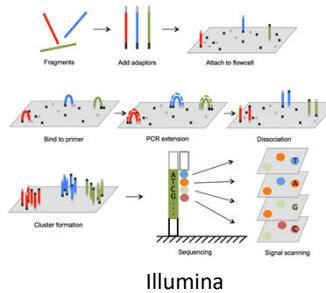
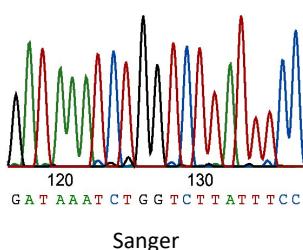


Sample types

1. Biological replicates
2. Technical replicates
3. Extraction negatives
4. PCR negatives

PCR-induced errors

- **PCR mutations:** polymerase enzymes introduce erroneous nucleotides now and then, even those enzymes with proof-reading activity
 - Dependent on the technology whether these becomes «visible»
 - In classic (direct) Sanger sequencing such errors become «diluted»
 - In methods where your final sequences are derived from one single DNA template, they become visible and must be corrected for!



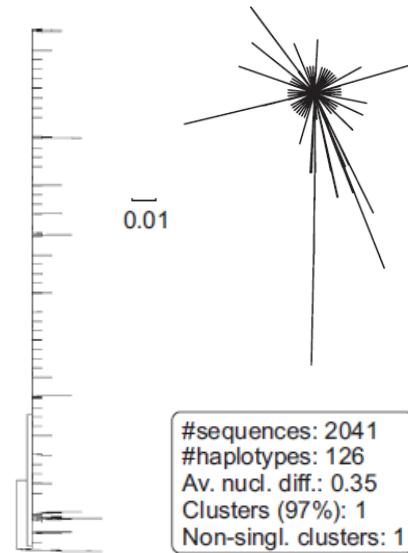
PCR-induced errors

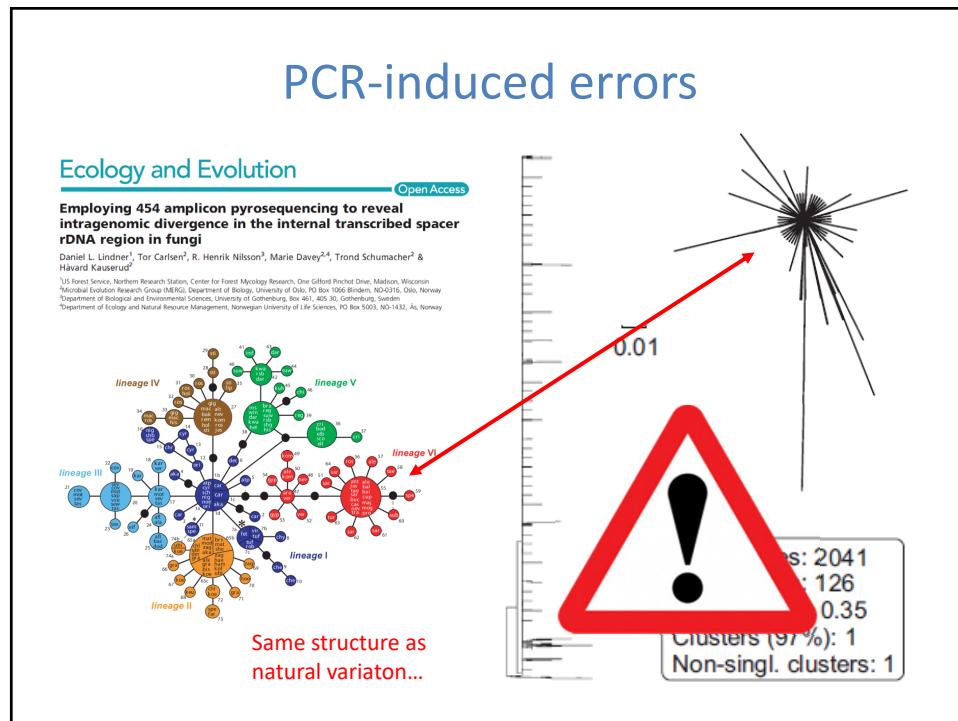
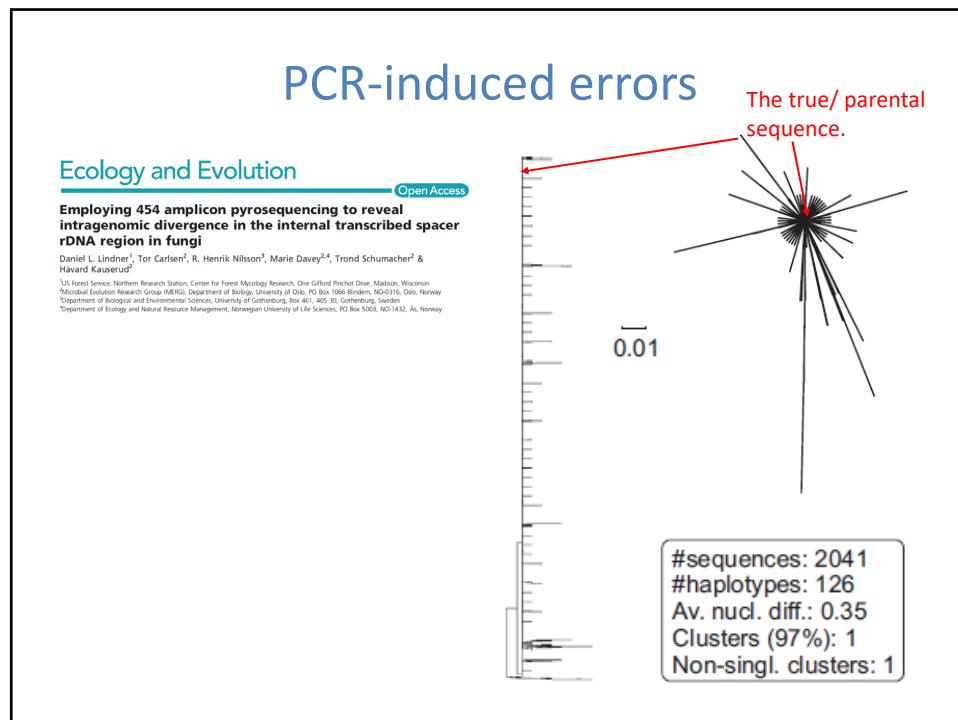
Ecology and Evolution Open Access

Employing 454 amplicon pyrosequencing to reveal intragenomic divergence in the internal transcribed spacer rDNA region in fungi

Daniel L. Lindner¹, Tor Carter², R. Henrik Nilsson³, Marie Davey^{2,4}, Trond Schumacher² & Howard K. Burge¹

¹US Forest Service, Northern Research Station, Center for Forest Mycology Research, One Gifford Pinchot Drive, Madison, Wisconsin
²Microbial Evolution Research Group (MERG), Department of Biology, University of Oslo, PO Box 1066 Blindern, NO-0316, Oslo, Norway
³Department of Biological and Environmental Sciences, University of Gothenburg, Box 460, 402 30, Gothenburg, Sweden
⁴Department of Ecology and Natural Resource Management, Norwegian University of Life Sciences, PO Box 5300, NO-1432, Ås, Norway





PCR-induced errors

- Chimeric sequences

MOLECULAR ECOLOGY
 RESOURCES
 Molecular Ecology Resources (2016)
 doi: 10.1111/1755-0998.12622

ITS all right mama: investigating the formation of chimeric sequences in the ITS2 region by DNA metabarcoding analyses of fungal mock communities of different complexities
 ANDERS BJØRNSGAARD A AS, MARIE LOUISE DAVEY and HÅVARD KAUSERUD
 Section for Genetics and Evolutionary Biology (Engen), Department of Biosciences, University of Oslo, P.O. Box 1066 Blindern, NO-0316 Oslo, Norway

The level of chimeric sequences depends on how variable the marker is!

Can reduce the problem with certain PCR settings, including long extension time and low number of cycles.

Polymerase enzyme

The Journal of Microbiology (2012) Vol. 50, No. 6, pp. 1071–1074
 Copyright © 2012, The Microbiological Society of Korea
 DOI 10.1007/s12275-012-3642-4

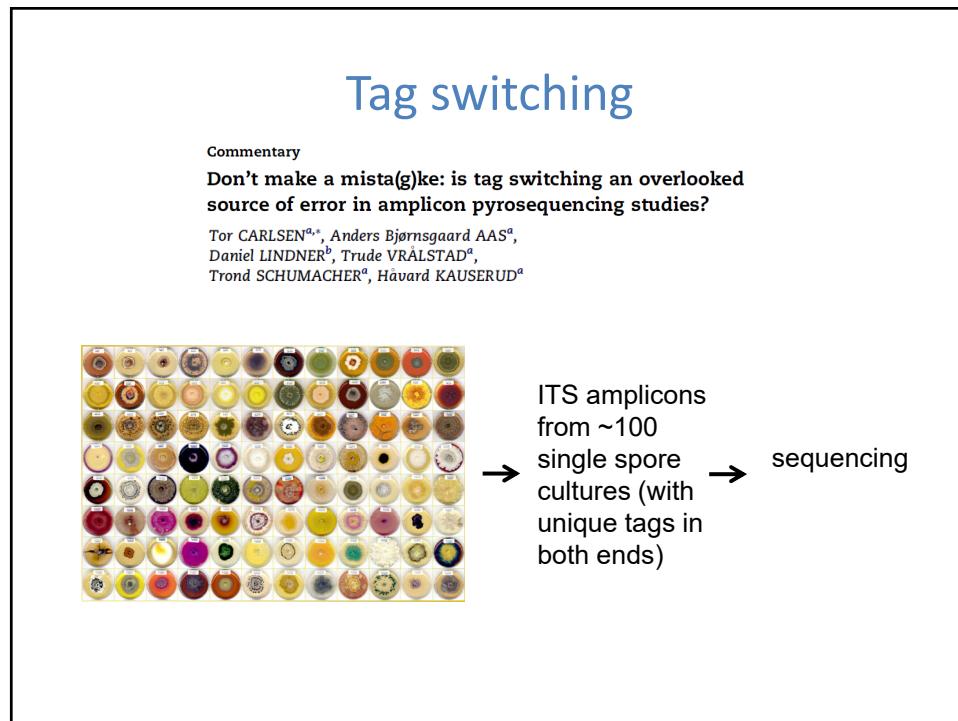
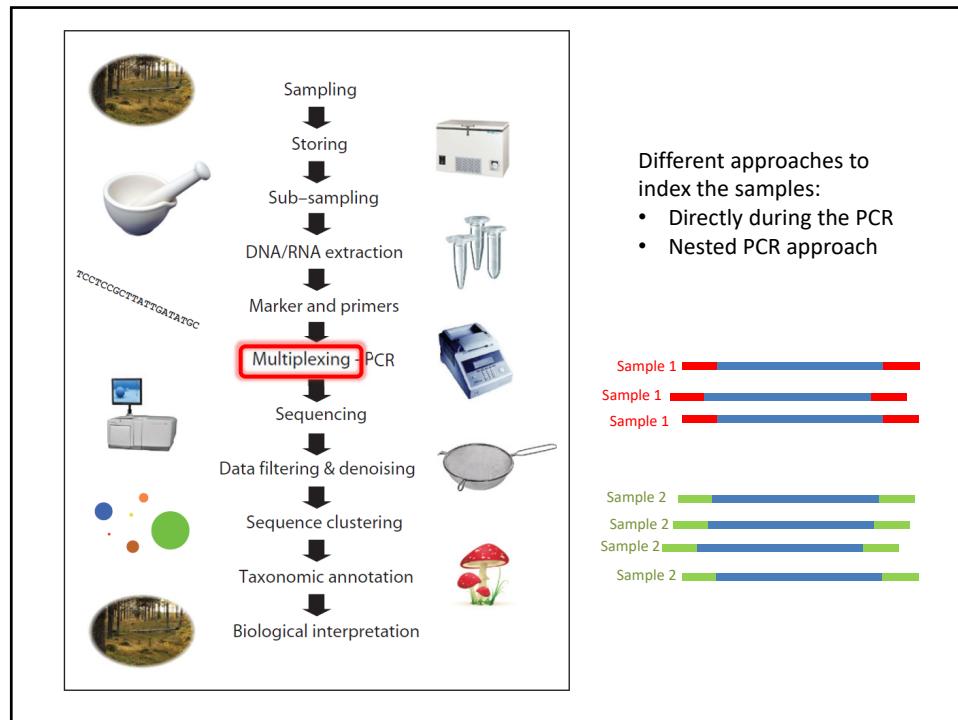
NOTE

Effects of PCR Cycle Number and DNA Polymerase Type on the 16S rRNA Gene Pyrosequencing Analysis of Bacterial Communities^a

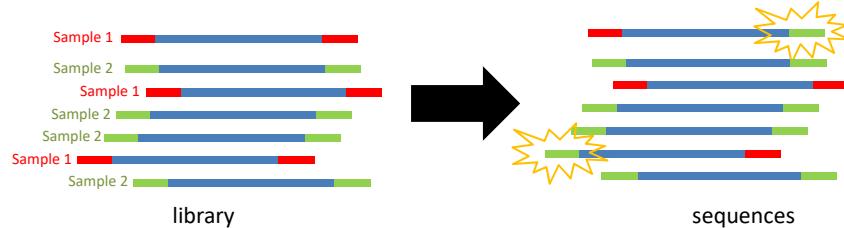
Jae Hyung Ahn, Byung-Yong Kim,
 Jackyeong Song, and Hang-Yeon Weon*

Agricultural Microbiology Division, National Academy of Agricultural Science, Rural Development Administration, Suwon 441-707, Republic of Korea
 (Received November 19, 2012 / Accepted December 11, 2012)

The bacterial richness was overestimated at increased PCR cycle number mostly due to the occurrence of chimeric sequences, and this was more serious with a DNA polymerase having proofreading activity than with *Taq* DNA polymerase. These results suggest that PCR cycle number must be kept as low as possible for accurate estimation of bacterial richness and that particular care must be taken when a DNA polymerase having proofreading activity is used.



Tag switching (tag jumping, bleeding, leaking, etc.)



Samples							
	1	2	3	4	5	6	7
OTU1	0	0	0	0	0	0	0
OTU2	2	0	10000	0	0	5	0
OTU3	0	0	0	0	0	0	0
OTU4	0	0	0	0	0	0	0
OTU5	0	0	0	0	0	0	0
OTU6	0	500	0	0	0	4	0
OTU7	0	0	0	0	0	0	0
OTU8	0	0	0	0	0	0	0
OTU9	0	0	23	0	0	80000	0
OTU10	0	0	0	0	0	0	0

Can lead to
numerous false
positives!



Tag switching

MOLECULAR ECOLOGY
RESOURCES

Tag jumps illuminated – reducing sequence-to-sample misidentifications in metabarcoding studies

IDA BÆR HOLM SCHNELL,^{*†} KRISTINE BOHMANN^{*‡} and M. THOMAS P. GILBERT^{*§}
^{*Centre for GeoGenetics, Natural History Museum of Denmark, University of Copenhagen, 1350 Copenhagen K, Denmark, †Center for Zoo and Wild Animal Health, Copenhagen Zoo, 1300 Frørslevsgade, Denmark, ‡School of Biological Sciences, University of Bristol, Bristol BS8 1UG, UK, §Trace and Environmental DNA Laboratory, Department of Environment and Agriculture, Curtin University, Perth, Western Australia 6102, Australia}

"We found that an average of 2.6% and 2.1% of sequences had tag combinations, which could be explained by tag jumping..."

Accurate multiplexing and filtering for high-throughput amplicon sequencing

High-throughput amplicon sequencing

Esling Philippe^{1,2,*}, Lejzerowicz Franck¹ and Pawlowski Jan¹

Received April 22, 2014; Revised January 29, 2015; Accepted January 29, 2015

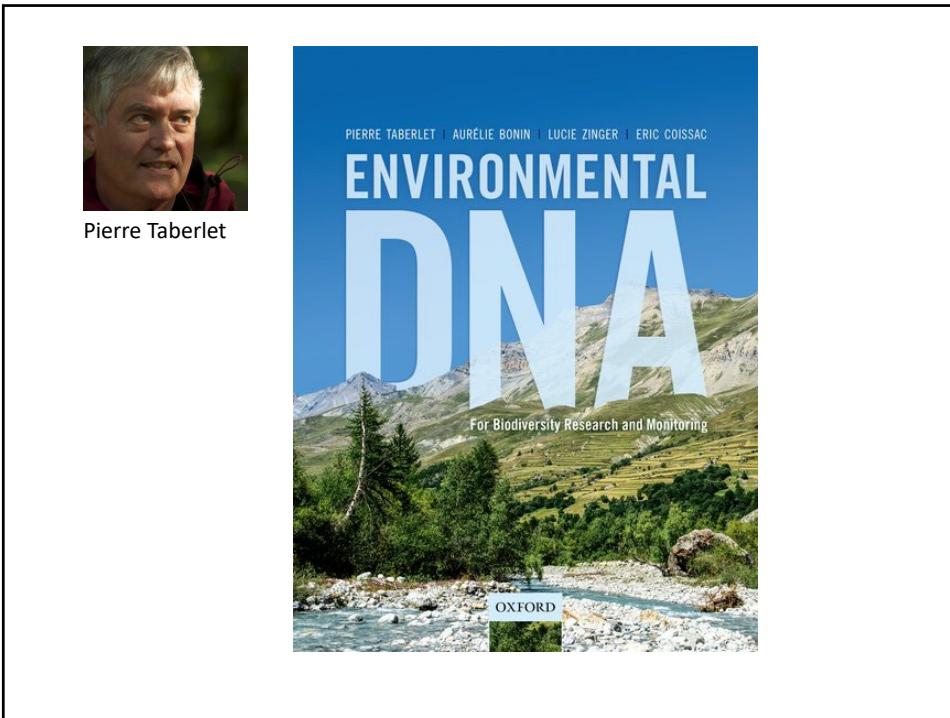
Up to 28.2% of the unique sequences correspond to undetectable (critical) mistakes in single- or saturated double-tagging libraries.

Tag switching

- The problem can be reduced or controlled for by:
 - Tagging in both ends with unique tag combinations
 - Rinse the PCR amplicons thoroughly
 - Include positive controls during PCR (mock community) → can better identify the level of switching/leakage
 - Avoid PCR step during the final library preparations steps before sequencing (i.e. when adaptors are introduced)

Tag switching

- The problem can be reduced or controlled for by:
 - Tagging in both ends with unique tag combinations
 - Rinse the PCR amplicons thoroughly
 - **Avoid PCR step during the final library preparations steps before sequencing (i.e. when adaptors are introduced)**
 - Include positive controls during PCR (mock community) → can better identify the level of switching/leakage



Tag switching

- The problem can be reduced or controlled for by:
 - Tagging in both ends with unique tag combinations
 - Rinse the PAvoid PCR step during the final library preparations steps before sequencing (i.e. when adaptors are introduced)
 - CR amplicons thoroughly
 - **Include positive controls during PCR (mock community) → can better identify the level of switching/leakage**

ZymoBIOMICS Microbial Community Standards

To improve the quality and reproducibility of metagenomic analyses, Zymo Research has endeavored to develop microbial reference materials. The ZymoBIOMICS Microbial Community Standard is the first commercially available standard for microbiomics and metagenomics research. This standard is composed of a complex mixture of pure bacterial and yeast strains with varying growth rates, Gram-negative and Gram-positive bacteria and yeast with varying sizes and cells with different wall properties. The wide range of organisms with different properties enables characterization, optimization, and validation of lysis methods such as bead beating. It can be used as a reference material for the development of new sequencing methods, as well as for the validation of existing methods. This standard is optimized and validated. A mock microbial DNA community standard allows researchers to focus the optimization after the step of DNA extraction.

Catalog #	Product	Size
DK300	ZymoBIOMICS Microbial Community Standard	10 Picograms
DK305	ZymoBIOMICS Microbial Community DNA Standard	200 ng
DK310	ZymoBIOMICS Microbial Community DNA Standard II [Log Distribution]	10 Picograms
DK311	ZymoBIOMICS Microbial Community DNA Standard II [Log Distribution]	200ng/20μl
DK320	ZymoBIOMICS Spike-in Control I (High Microbial Load)	25 Picograms
DK321	ZymoBIOMICS Spike-in Control I (Low Microbial Load)	25 Picograms
DK322	ZymoBIOMICS HhW DNA Standard	5000 ng
DK323	ZymoBIOMICS Fecal Reference with TruBact™ Technology	10 picograms
PAK00	ZymoBIOMICS Cut Microbiome Standard	10 picograms

Sample types

1. Biological replicates
2. Technical replicates
3. Extraction negatives
4. PCR negatives
5. Positive control (mock community)

EDITORIAL

MOLECULAR ECOLOGY WILEY

DNA metabarcoding—Need for robust experimental designs to draw sound ecological conclusions

Zinger et al. 2019. Molecular Ecology Resources

Sample types

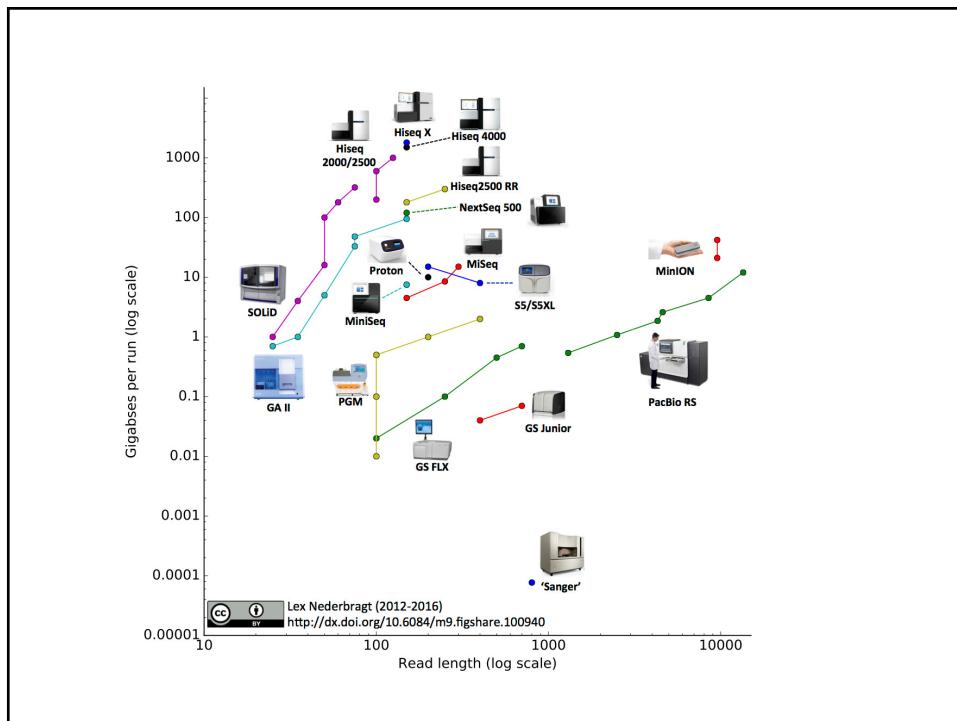
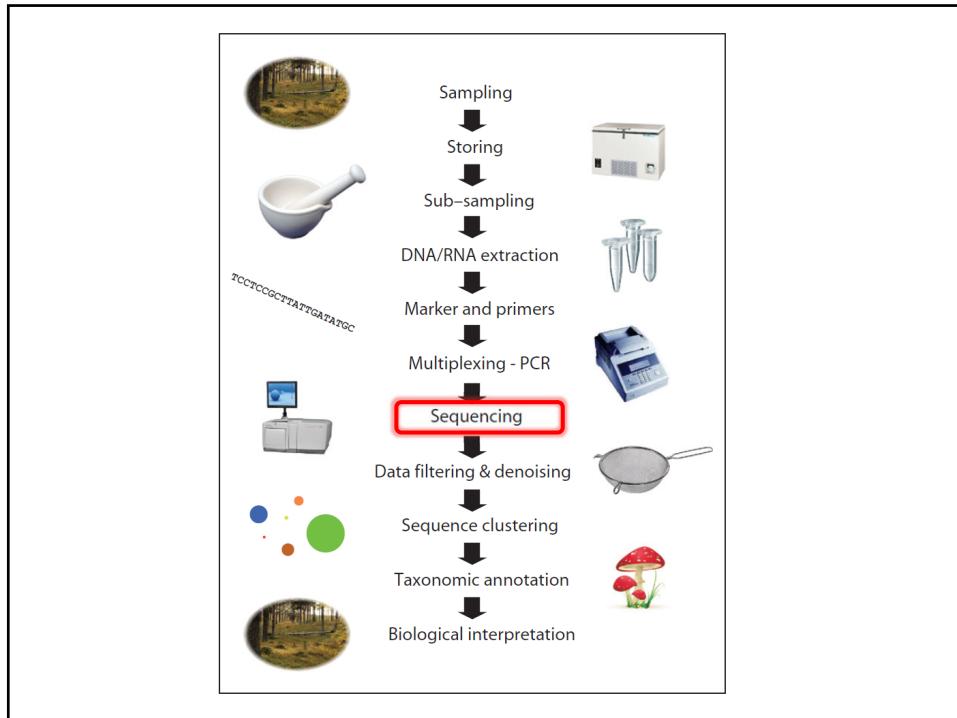
1. Biological replicates
2. Technical replicates → Not strictly required
3. Extraction negatives | → Only sequence if positive amplicon?
4. PCR negatives |
5. Positive control (mock community)

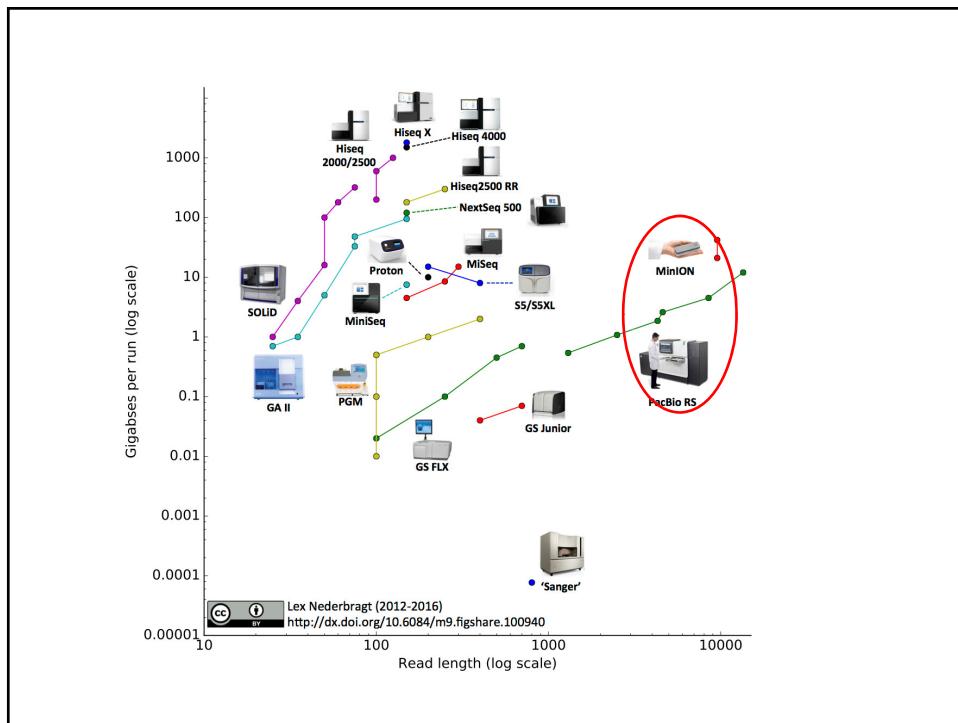
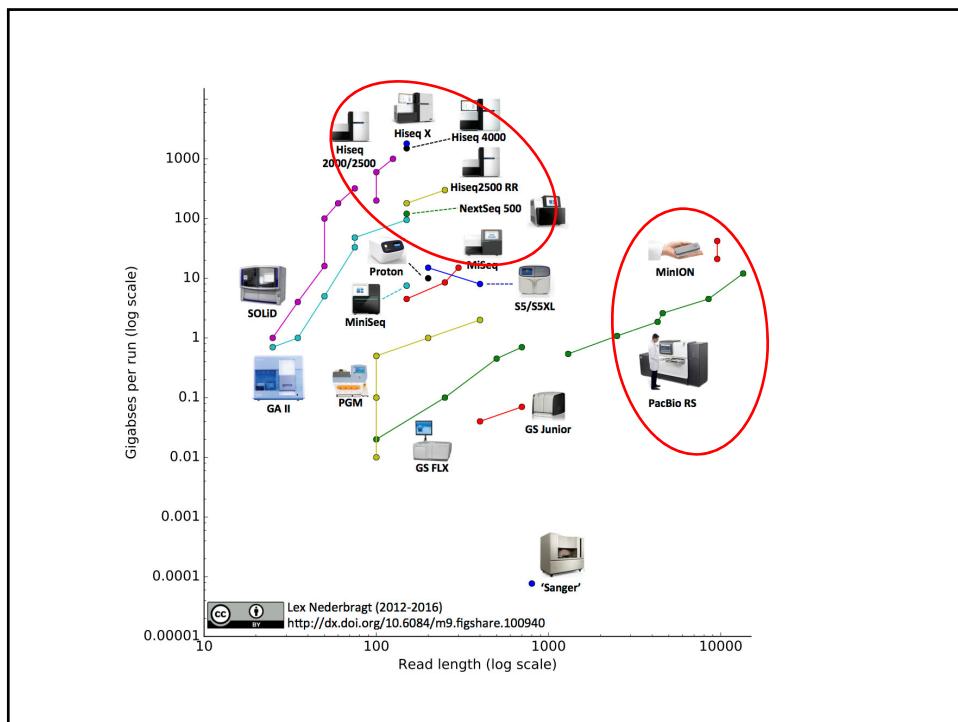
EDITORIAL

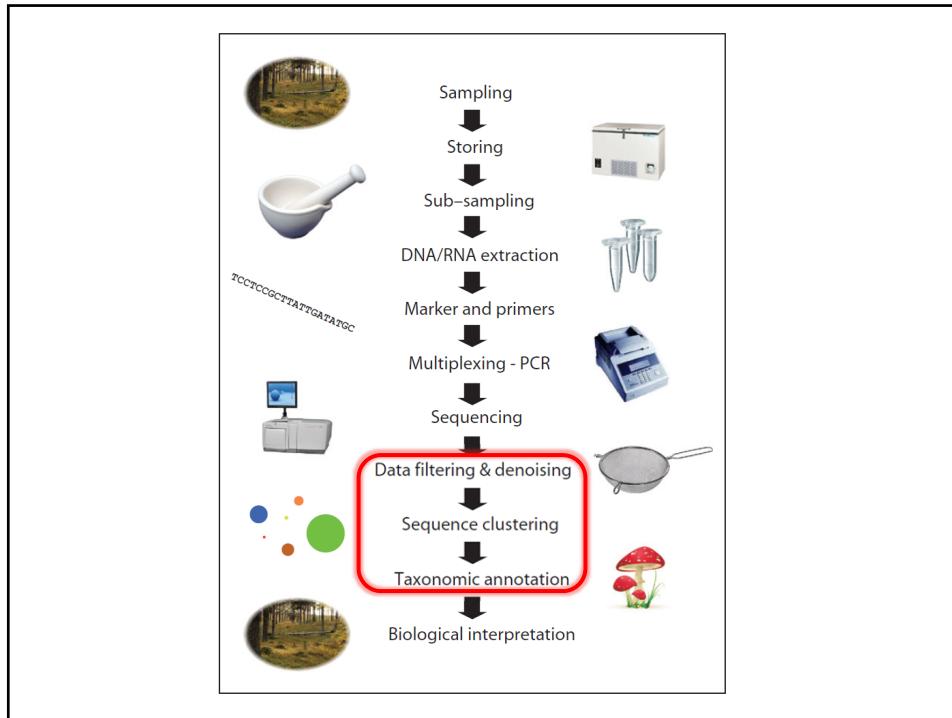
MOLECULAR ECOLOGY WILEY

DNA metabarcoding—Need for robust experimental designs to draw sound ecological conclusions

Zinger et al. 2019. Molecular Ecology Resources



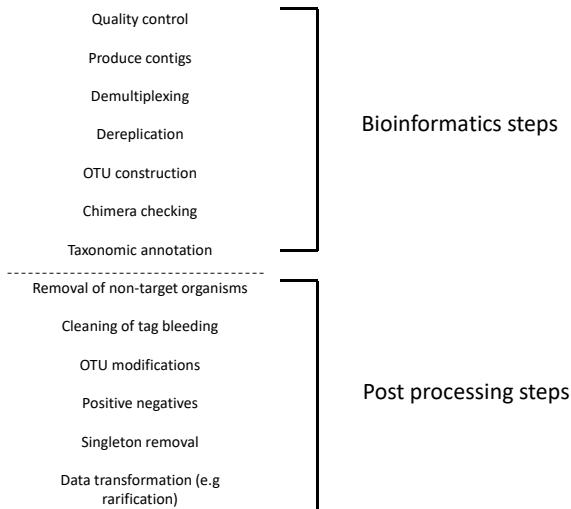




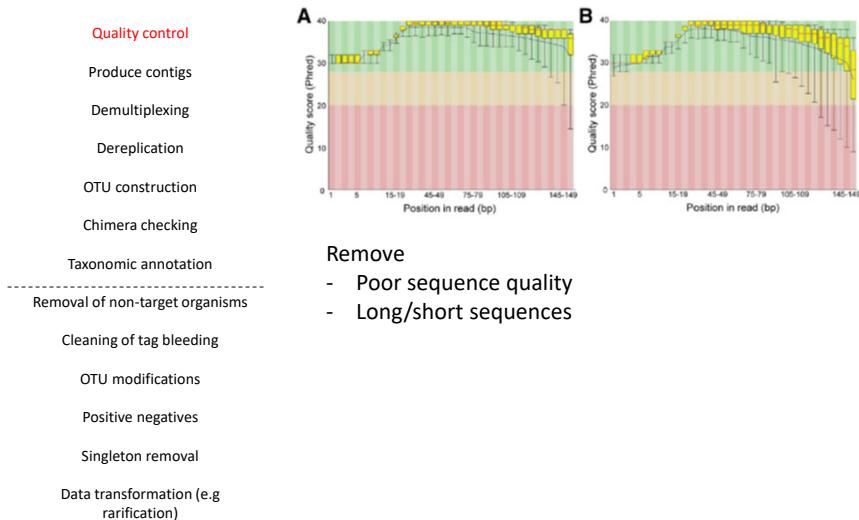
Bioinformatics – main steps

- (The order of steps depends somewhat on the pipeline/programs)
- Quality control
 - Produce contigs
 - Demultiplexing
 - Dereplication
 - OTU construction
 - Chimera checking
 - Taxonomic annotation
-
- Removal of non-target organisms
- Cleaning of tag bleeding
 - OTU modifications
 - Positive negatives
 - Singleton removal
 - Data transformation (e.g. rarification)

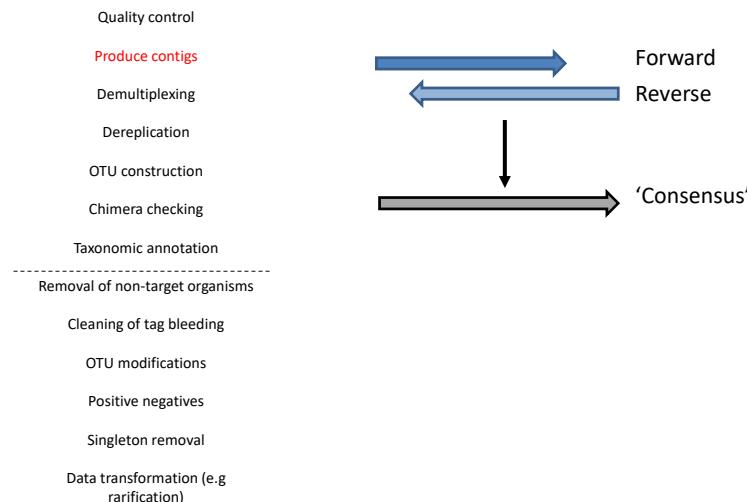
Bioinformatics – main steps



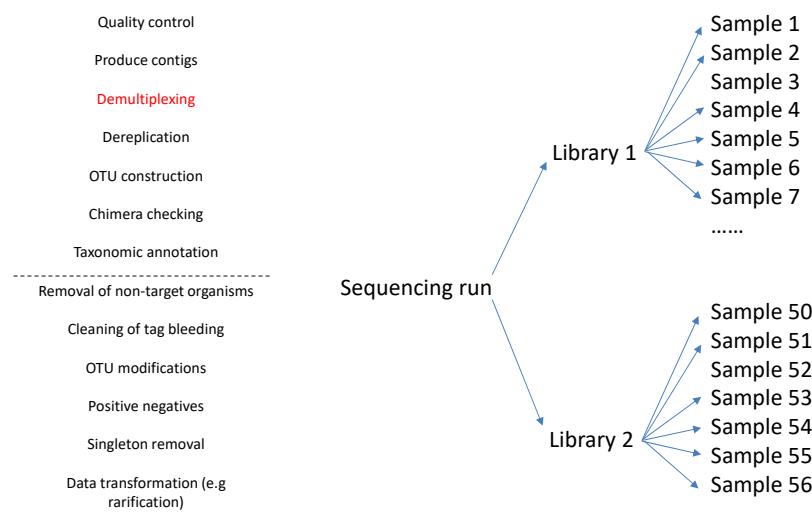
Bioinformatics – main steps



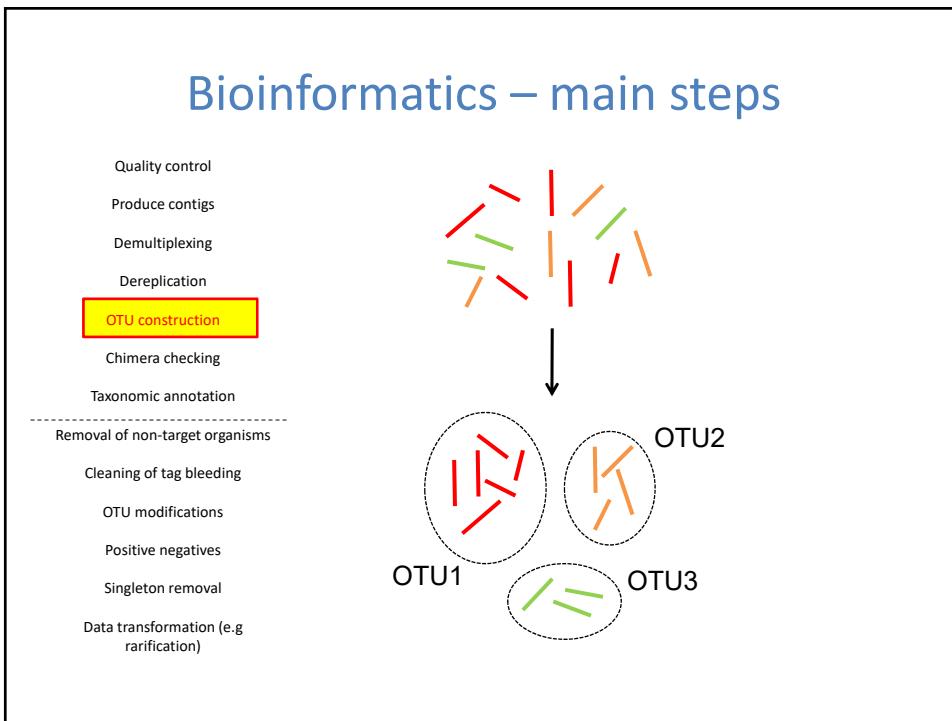
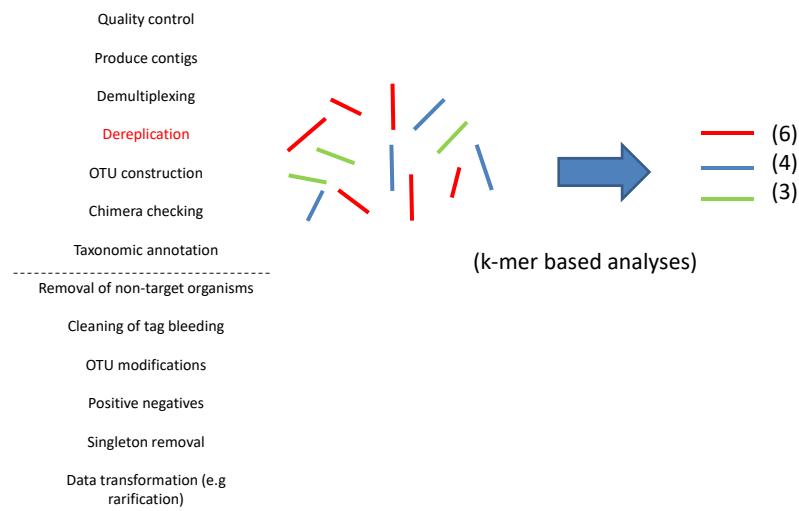
Bioinformatics – main steps

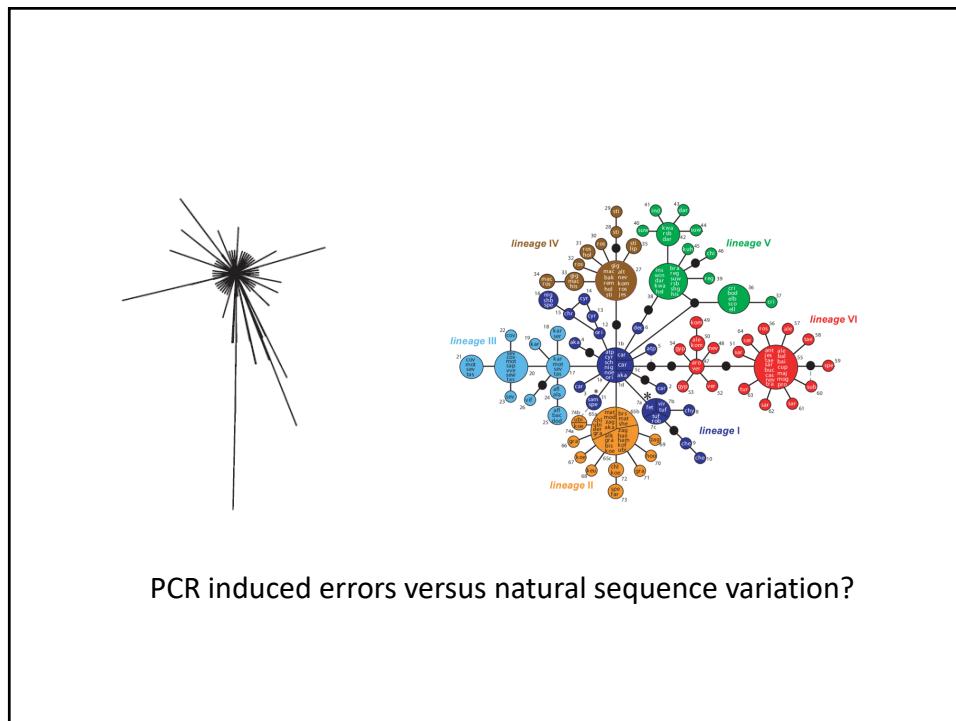
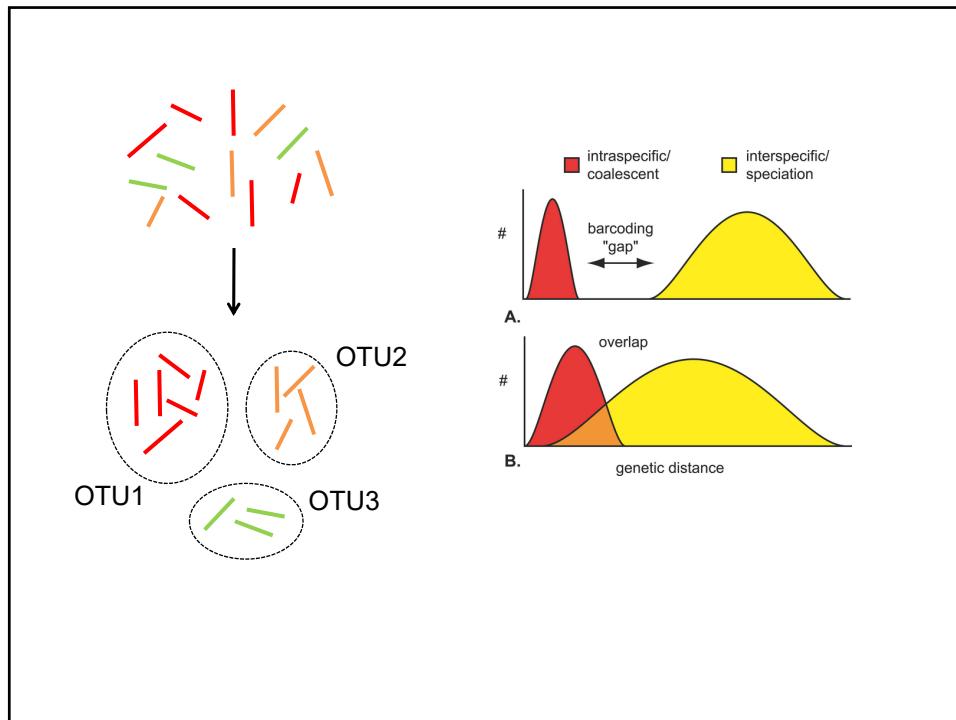


Bioinformatics – main steps



Bioinformatics – main steps





Bioinformatics – main steps

Quality control
Produce contigs
Demultiplexing
Dereplication
OTU construction
Chimera checking
Taxonomic annotation

Removal of non-target organisms
Cleaning of tag bleeding
OTU modifications
Positive negatives
Singleton removal
Data transformation (e.g. rarification)

de novo versus closed (reference based) OTU construction?

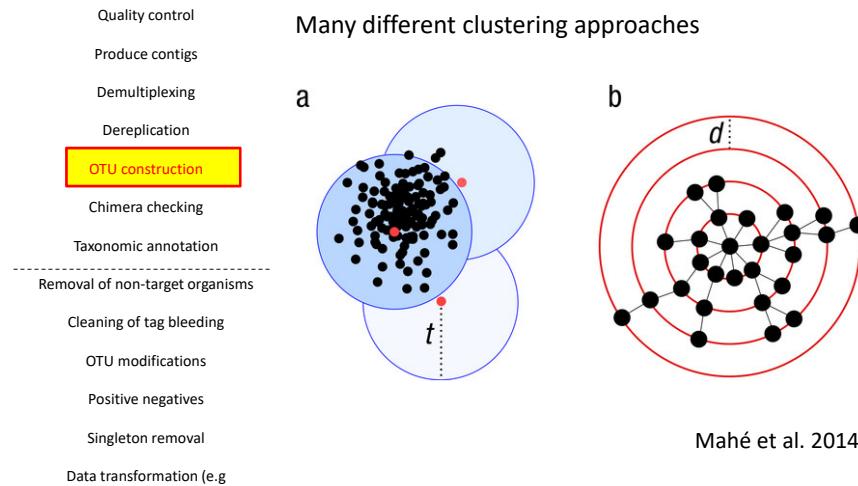
Bioinformatics – main steps

Quality control
Produce contigs
Demultiplexing
Dereplication
OTU construction
Chimera checking
Taxonomic annotation

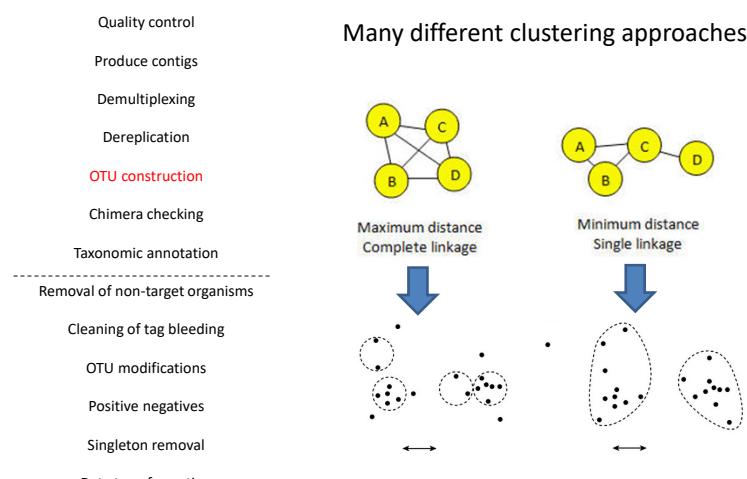
Removal of non-target organisms
Cleaning of tag bleeding
OTU modifications
Positive negatives
Singleton removal
Data transformation (e.g. rarification)

de novo versus closed (reference based) OTU construction?

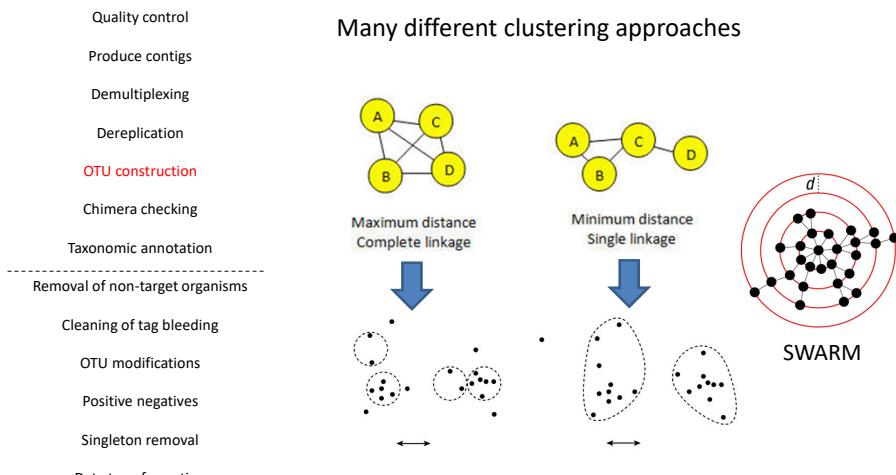
Bioinformatics – main steps



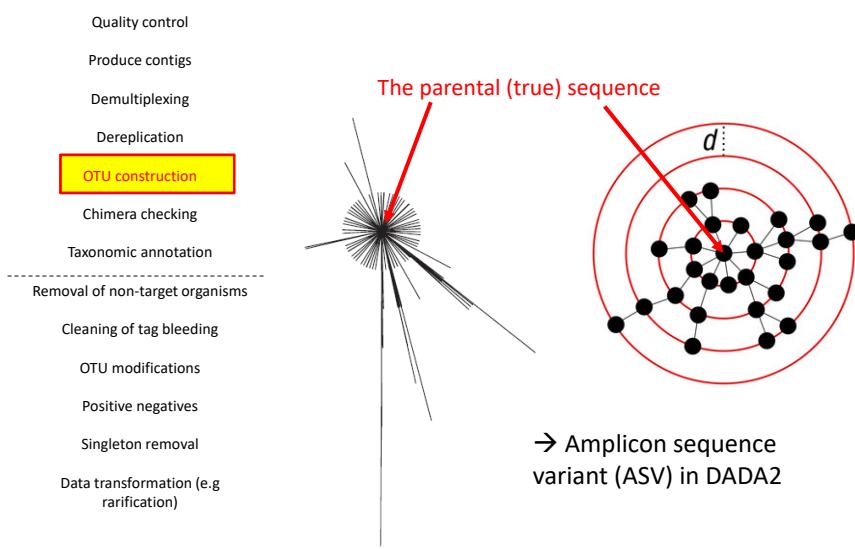
Bioinformatics – main steps



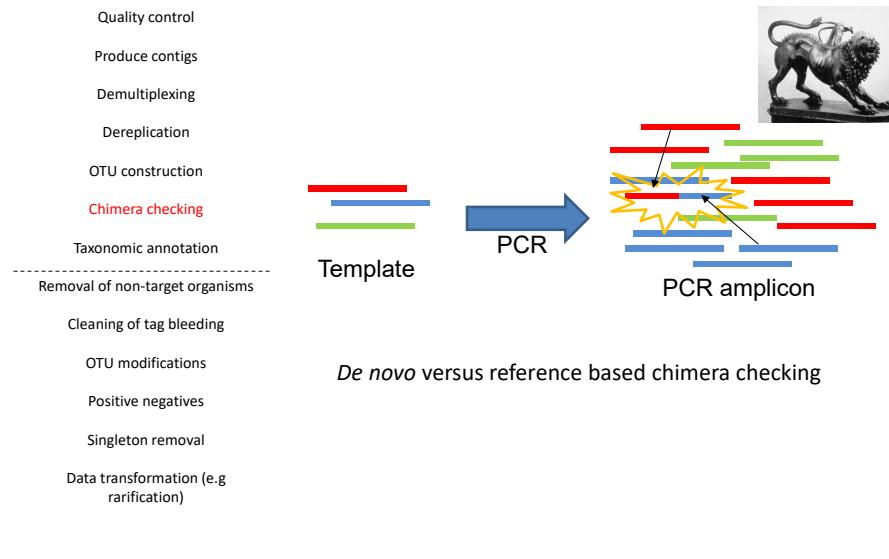
Bioinformatics – main steps



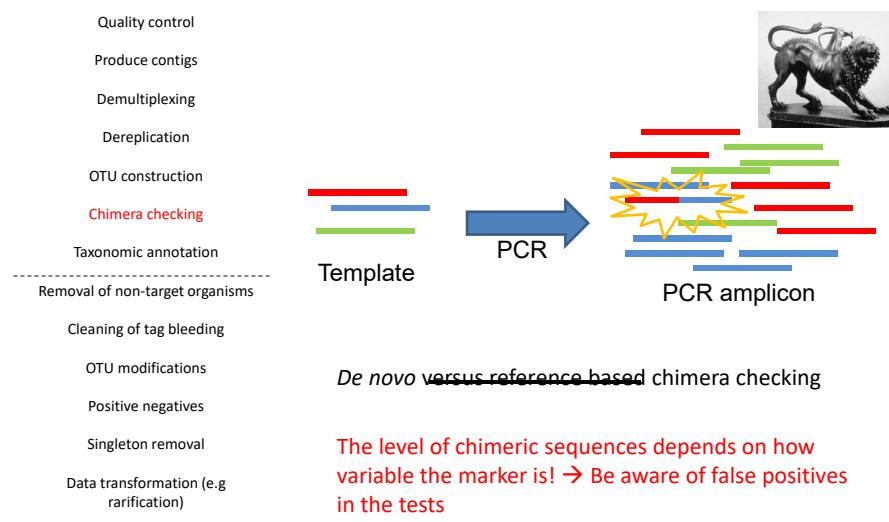
Bioinformatics – main steps



Bioinformatics – main steps



Bioinformatics – main steps



Bioinformatics – main steps

Quality control

Produce contigs

Demultiplexing

Dereplication

OTU construction

Chimera checking

Taxonomic annotation

Removal of non-target organisms

Cleaning of tag bleeding

OTU modifications

Positive negatives

Singleton removal

Data transformation (e.g.
rarification)



Simple matching (blast) → probabilistic
assignment (protax)

Bioinformatics – main steps

Quality control

Produce contigs

Demultiplexing

Dereplication

OTU construction

Chimera checking

Taxonomic annotation

Removal of non-target organisms

Cleaning of tag bleeding

OTU modifications

Positive negatives

Singleton removal

Data transformation (e.g.
rarification)

FROM THE COVER

MOLECULAR ECOLOGY
RESOURCES WILEY

Assessment of current taxonomic assignment strategies for
metabarcoding eukaryotes

Jose S. Hleep^{1,2,3} | Joanne E. Littlefair^{1,4} | Dirk Steinke⁵ | Paul D. N. Hebert⁵ |
Melania E. Cristescu¹

Hleep et al. 2021. Mol Ecol Resources

Bioinformatics – main steps

Quality control

Produce contigs

Demultiplexing

Dereplication

OTU construction

Chimera checking

Taxonomic annotation

Removal of non-target organisms

Cleaning of tag bleeding

OTU modifications

Positive negatives

Singleton removal

Data transformation (e.g.
rariﬁcation)

Bioinformatics – main steps

Quality control



Produce contigs



Demultiplexing



Dereplication



OTU construction



Chimera checking

Taxonomic annotation

Removal of non-target organisms

Cleaning of tag bleeding

Samples

	1	2	3	4	5	6	7
OTU1	0	0	0	0	0	0	0
OTU2	2	0	10000	0	0	5	0
OTU3	0	0	0	0	0	0	0
OTU4	0	0	0	0	0	0	0
OTU5	0	0	0	0	0	0	0
OTU6	0	500	0	0	0	4	0
OTU7	0	0	0	0	0	0	0
OTU8	0	0	0	0	0	0	0
OTU9	0	0	23	0	0	30000	0
OTU10	0	0	0	0	0	0	0

OTU modifications

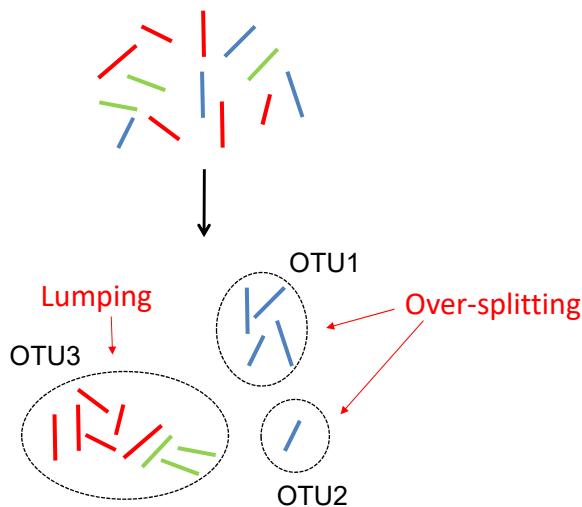
Positive negatives

Singleton removal

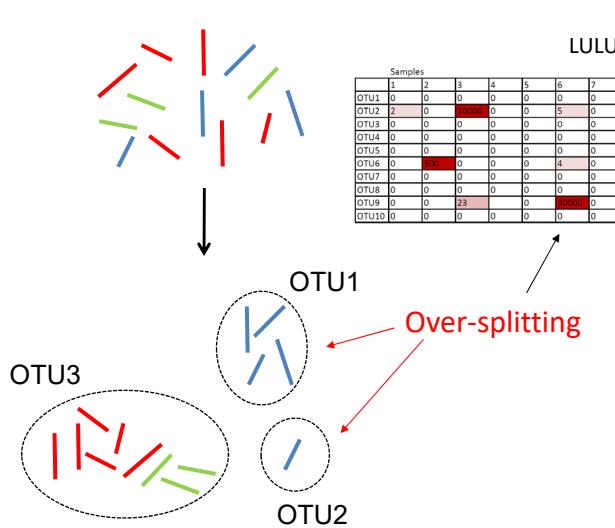
Data transformation

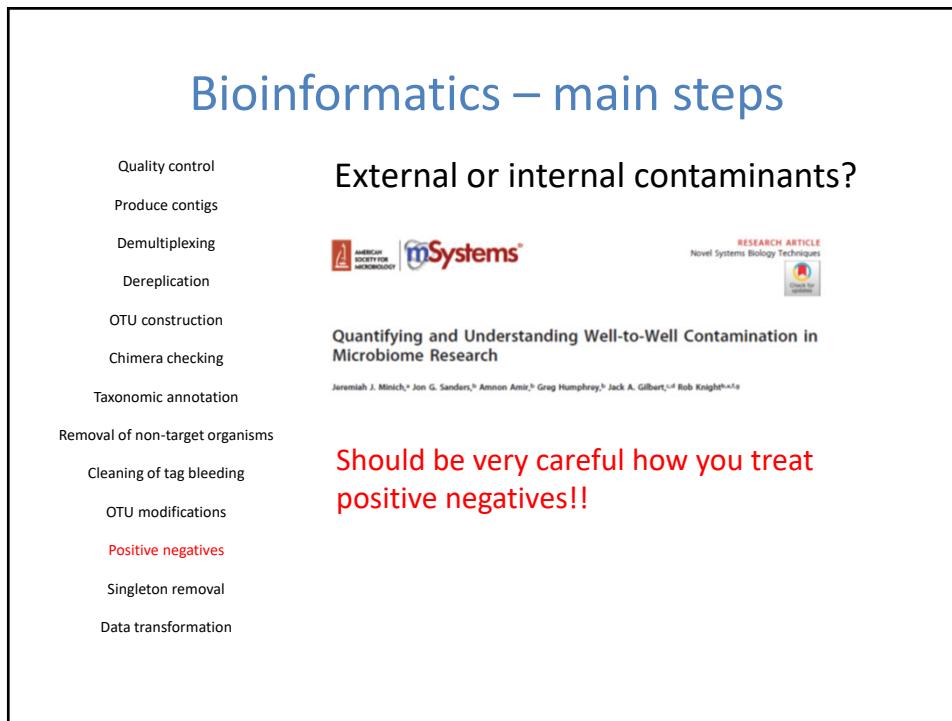
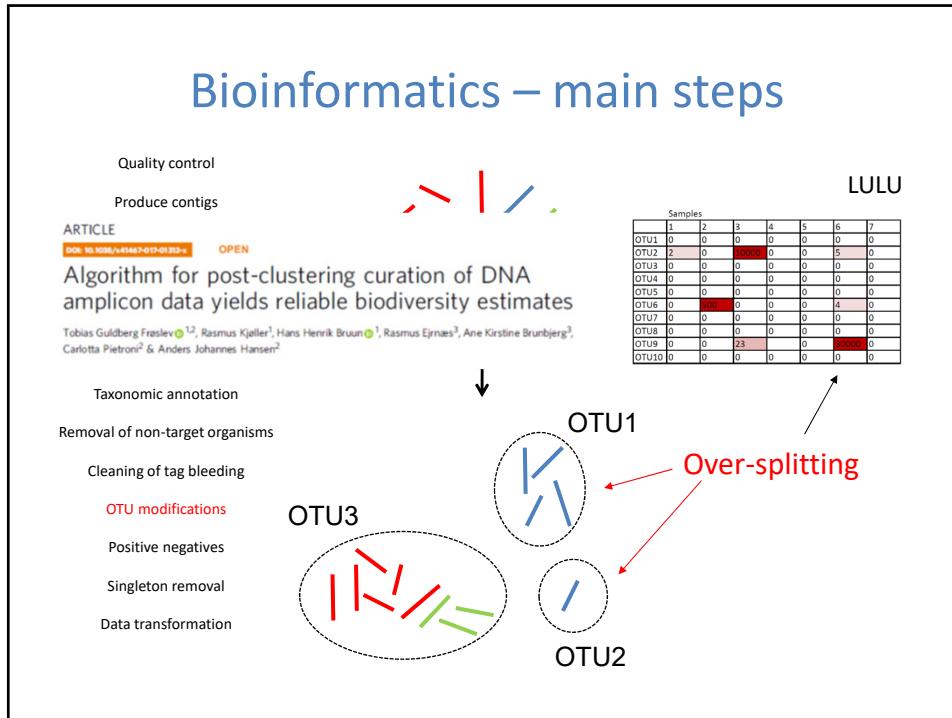
Bioinformatics – main steps

Quality control
 Produce contigs
 Demultiplexing
 Dereplication
 OTU construction
 Chimera checking
 Taxonomic annotation
 Removal of non-target organisms
 Cleaning of tag bleeding
OTU modifications
 Positive negatives
 Singleton removal
 Data transformation



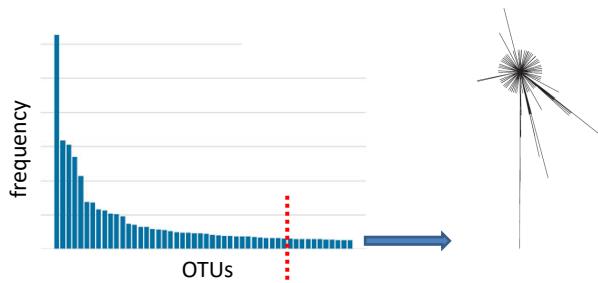
Quality control
 Produce contigs
 Demultiplexing
 Dereplication
 OTU construction
 Chimera checking
 Taxonomic annotation
 Removal of non-target organisms
 Cleaning of tag bleeding
OTU modifications
 Positive negatives
 Singleton removal
 Data transformation





Bioinformatics – main steps

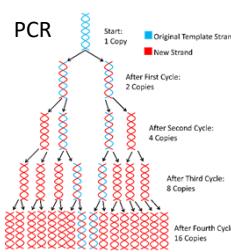
Quality control
 Produce contigs
 Demultiplexing
 Dereplication
 OTU construction
 Chimera checking
 Taxonomic annotation
 Removal of non-target organisms
 Cleaning of tag bleeding
 OTU modifications
 Positive negatives
Singleton removal
 Data transformation



What is a ‘singleton’? → Depends on your sequencing depth and quality of your data. Should also take study aim into consideration

Quality control
 Produce contigs
 Demultiplexing
 Dereplication
 OTU construction
 Chimera checking
 Taxonomic annotation
 Removal of non-target organisms
 Cleaning of tag bleeding
 OTU modifications
 Positive negatives
Singleton removal
 Data transformation

	Samples	1	2	3	4	5	6	7
OTU1	0	0	0	0	0	0	0	0
OTU2	2	0	0	10000	0	0	5	0
OTU3	0	0	0	0	0	0	0	0
OTU4	0	0	0	0	0	0	0	0
OTU5	0	0	0	0	0	0	0	0
OTU6	0	500	0	0	0	0	4	0
OTU7	0	0	0	0	0	0	0	0
OTU8	0	0	0	0	0	0	0	0
OTU9	0	0	23	0	0	0	80000	0
OTU10	0	0	0	0	0	0	0	0



Be careful with resampling and transformations!

Check the effect from various data treatments options on the results!

Depends on the study aims!

The importance of controls

1. Biological replicates
2. Technical replicates
3. Extraction negatives
4. PCR negatives
5. Positive control (mock community)



Different purposes

Conclusions

- Which methods to use? → No general answer – it is context dependent. You must argue for your choices!

