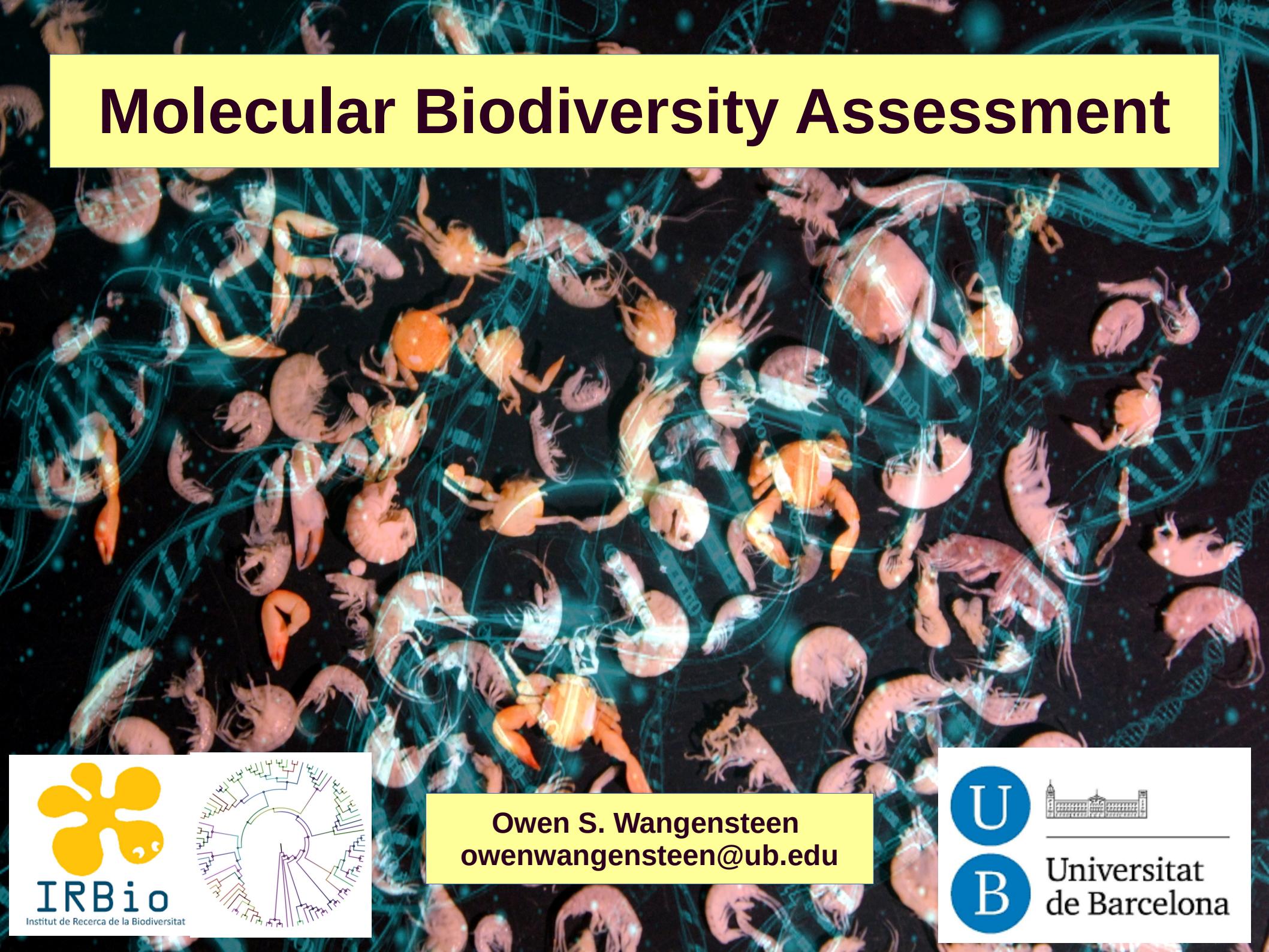


Molecular Biodiversity Assessment

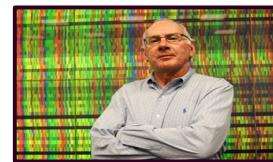


OUTLINE:

1. Importance of biodiversity assessment and biomonitoring



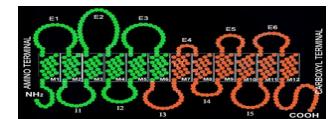
2. Molecular approaches to biodiversity: DNA-barcoding, DNA-metabarcoding. Other metagenomics techniques.



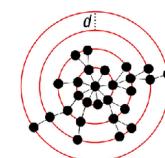
3. The metabarcoding approach: Community DNA vs extra-organismal DNA. Quantitative value. Metabarcoding markers, metabarcoding primer design.



4. Metabarcoding workflow and pipelines



4.1. Sampling and pre-processing



4.2. DNA extraction and PCR

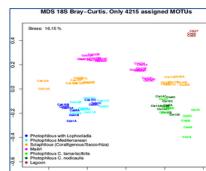


4.3. Bioinformatics: demultiplexing and QC

4.4. Bioinformatics: MOTU definition / denoising /clustering

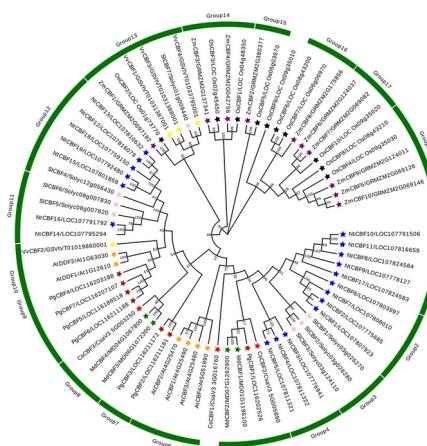
4.5. Bioinformatics: taxonomic assignment

4.6. Bioinformatics: final dataset refining



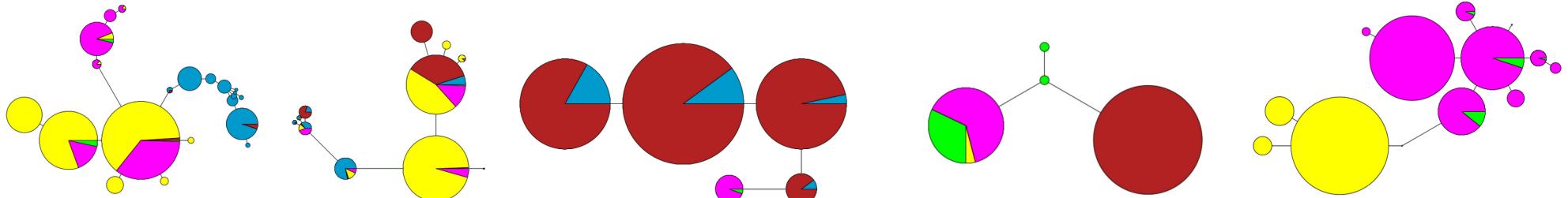
1. IMPORTANCE OF BIODIVERSITY ASSESSMENT AND BIOMONITORING

What is a class on **Biodiversity Assessment** like you doing in a course on **Phylogenomics and Population Genomics** like this?



1- Classical (morphological) biodiversity assessment methods are on their way to be replaced by **molecular biodiversity assessment** methods. Bioinformatic analyses of HT-sequencing datasets and taxonomic assignment of detected sequences is crucial for these applications.

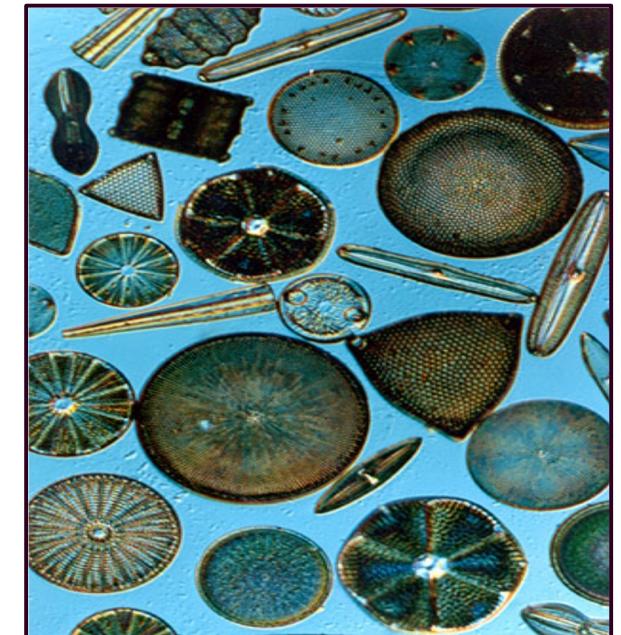
2- Some molecular biodiversity assessment methods can be used to infer population genetic patterns for many species at one time, even for species without taxonomic assignment (**metaphylogeography**)



BIOMONITORING AND ECOLOGICAL QUALITY

Some examples:

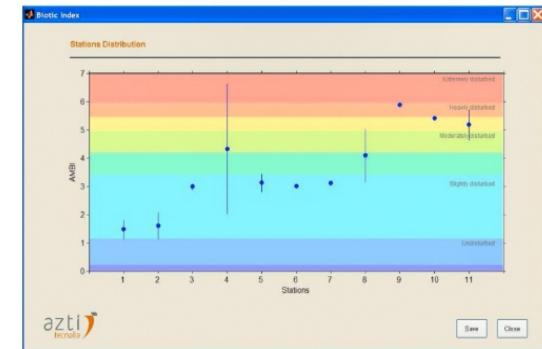
- River benthic macroinvertebrates
- Marine benthic macroinvertebrates
- Freshwater microalgae (diatoms)



ECOLOGICAL QUALITY INDICES BASED ON BIODIVERSITY

Some examples:

- RICT: River benthic macroinvertebrates
- AMBI: Marine benthic macroinvertebrates
- IBD: Freshwater diatoms



Drawbacks:

- Need of an experienced taxonomist
- Need of adaptive implementation to different biogeographical areas
- Slow! Thorough analysis of some complex samples can take weeks!
- Subjectivity, lack of repeatability and lack of traceability

ERRORS IN MORPHOLOGY-BASED ECOLOGICAL QUALITY INDICES

J. N. Am. Benthol. Soc., 2010, 29(4):1279–1291
© 2010 by The North American Bentholological Society
DOI: 10.1899/09-183.1
Published online: 7 September 2010

First audit of macroinvertebrate samples from an EU Water Framework Directive monitoring program: human error greatly lowers precision of assessment results

Peter Haase^{1,3}, Steffen U. Pauls^{1,2,4}, Karin Schindelhütte^{1,5}, AND
Andrea Sundermann^{1,6}

Errors and biases:

- 29 % of specimens and 21% of taxa were overlooked by the analysts
- 30 % of assigned taxa differed between analysts and auditors
- 34 % of audited samples were assigned to the wrong ecological quality class

In comparison, molecular biodiversity assessment methods:

- No need of an experienced taxonomist. A lab technician can do all the work
- Fast and scalable! Hundreds of samples can be multiplexed in the same run
- Objectivity, repeatability, perfect traceability, sequence datasets can be stored forever and they are accessible online for everyone
- Accuracy of the assessment depends on the existence of good reference sequence databases

DNA-BARCODING



Received 29 July 2002
Accepted 30 September 2002
Published online 8 January 2003

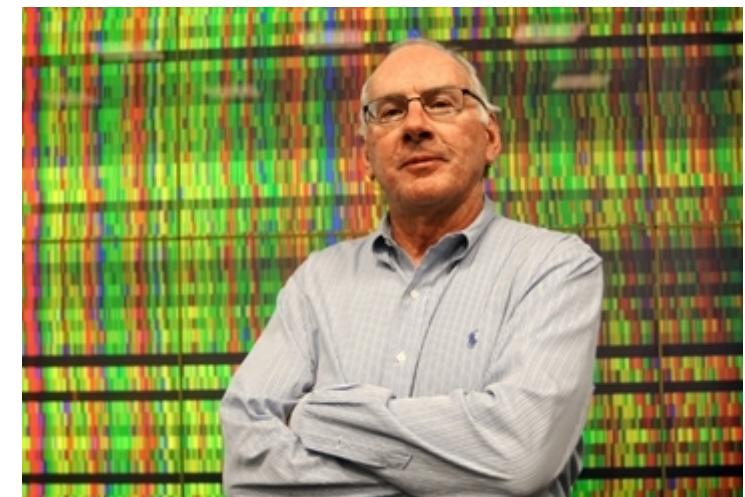
Biological identifications through DNA barcodes

Paul D. N. Hebert¹, Alina Cywinska, Shelley L. Ball
and Jeremy R. deWaard

Department of Zoology, University of Guelph, Guelph, Ontario N1G 2W1, Canada

Although much biological research depends upon species diagnoses, taxonomic expertise is collapsing. We are convinced that the sole prospect for a sustainable identification capability lies in the construction of systems that employ DNA sequences as taxon ‘barcodes’. We establish that the mitochondrial gene cytochrome *c* oxidase I (COI) can serve as the core of a global bioidentification system for animals. First, we demonstrate that COI profiles, derived from the low-density sampling of higher taxonomic categories, ordinarily assign newly analysed taxa to the appropriate phylum or order. Second, we demonstrate that species-level assignments can be obtained by creating comprehensive COI profiles. A model COI profile, based upon the analysis of a single individual from each of 200 closely allied species of lepidopterans, was 100% successful in correctly identifying subsequent specimens. When fully developed, a COI identification system will provide a reliable, cost-effective and accessible solution to the current problem of species identification. Its assembly will also generate important new insights into the diversification of life and the rules of molecular evolution.

Keywords: molecular taxonomy; mitochondrial DNA; animals; insects; sequence diversity; evolution



Paul Hebert (2003)

A short DNA fragment (~600 pb) of a highly variable genetic locus is enough to unequivocally identify a species with 100% success

THE BARCODE OF LIFE DATABASE

www.boldsystems.org

The screenshot shows the BOLD Systems homepage. At the top, there's a navigation bar with links for Databases, Taxonomy, Identification, Workbench, and Resources, along with a Log In button. Below the navigation is a world map where red dots represent barcode collection sites, appearing more densely in North America, Europe, and Asia. A callout box in the bottom right corner of the map area displays "Barcodes: 4,494,495" and "Per Site: 1000 | 100 | 10 | 1". Below the map is a search bar with a "Search" button. The main content area features four circular icons with accompanying text: "Public Data Portal" (data retrieval interface), "Barcode Index Numbers" (searchable database of BINs), "DNA Barcode Education Portal" (custom platform for educators and students), and "Workbench" (integrated data collection and analysis environment). Each icon has a small descriptive image next to it.

BOLDSYSTEMS Databases | Taxonomy | Identification | Workbench | Resources Log In

Barcodes: 4,494,495
Per Site: 1000 | 100 | 10 | 1

Taxonomy ▾ Search

Public Data Portal:
A data retrieval interface that allows for searching over 1.7M public records in BOLD using multiple search criteria including, but not limited to, geography, taxonomy, and depository.

Barcode Index Numbers:
A searchable database of Barcode Index Numbers (BINs), sequence clusters that closely approximate species.

DNA Barcode Education Portal:
A custom platform for educators and students to explore barcode data and contribute novel barcodes to the BOLD database.

Workbench:
An integrated data collection and analysis environment that securely supports the assembly and validation of DNA barcodes and ancillary sequences.

Barcodes: 17,552,000
BINs: 1,109,000

**Species with barcodes: 354,000
(formally described)**

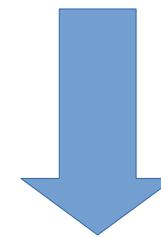
DNA-BARCODING IDENTIFICATION



Sanger
sequencing



```
AACATTATATTTATTTGGTATTCAGCAAGTATATTAGGAACCTCTTAGTTATTAAATTGAACTGAATT  
AGGAACACCGGGCTCATTAATTGGAGATGATCAAATTATAACTATTGTTACAGCTACGCTTTATTATAAT  
TTTTTTATAGTTACCTTATTATAATTGGAGGATTGATTAACTCCTTAACTTTAATTAGGAGCCCTGA  
CATAGCTTCCCACGAATAAAAATAAGATTGATTACCTCCTTCAACTCTTAACTGACATGGAAGAAG  
AATTGTAGAAAATGGAGCTGGAACCTGGTTGAACGTGTTACCCCCCTCTCTTCAACATTGACATGGAAGAAG  
ATCAGTAGATTAGTTATTTTCTTACATTAGCAGGTATCTCATCAATTAGGGCAATTAACTTCATCAC  
AACAAATTATTAATACGTTAAATAATATCATTGATCAAATACCTTATTGTTGAGCTGTTGGAATTAC  
AGCTTATTACTTCTTACCAAGTTAGCTGGAGCTATTACTATATTAAACAGATGAAATTAAA  
TACTTCTTTTGATCCTGCTGGAGGGGAGACCCTATTACCAACATTATT
```



REFERENCE DATABASE
www.boldsystems.org

[Go to public records in this BIN](#)

BIN DETAILS:

| | | | |
|-------------------|------------------|-------------------------------|----------------|
| BIN URI: | BOLD:AAA5810 | Average Distance: | 1.03% (p-dist) |
| DOI: | Pending | Maximum Distance: | 3.58% (p-dist) |
| Member Count: | 255 [212 Public] | Distance to Nearest Neighbor: | 1.77% (p-dist) |
| Barcode Compliant | 226 | | |
| Members: | | | |
| Founding Record: | | | |

NEAREST NEIGHBOR (NN) DETAILS:

| | | | |
|--------------------------|--|--------------------|----------------|
| Nearest BIN URI: | BOLD:ABZ2147 | Average Distance: | 0.12% (p-dist) |
| Member Count: | 2 | Maximum Distance: | 0.12% (p-dist) |
| Nearest Member: | GBGL7241-10 | Distance Variance: | 0% (p-dist) |
| Nearest Member Taxonomy: | Arthropoda, Insecta, Lepidoptera, Papilionidae, Papilioninae, Papilio, Papilio machaon | | |

TAXONOMY:

| | | |
|------------|---------------------------------|---|
| Phylum: | Arthropoda [255] | Q |
| Class: | Insecta [255] | Q |
| Order: | Lepidoptera [255] | Q |
| Family: | Papilionidae [255] | Q |
| Subfamily: | Papilioninae [255] | Q |
| Genus: | Papilio [255] | Q |
| Species: | Papilio machaon [226] | Q |
| | Papilio brevicauda [11] | Q |
| | Papilio saharae [8] | Q |
| | Papilio kahli [2] | Q |
| | Papilio machaon hippocrates [2] | Q |
| | Papilio machaon rinpoche [2] | Q |
| | Papilio hospiton [1] | Q |
| | Papilio sp. [1] | Q |
| | Papilio machaon syriacus [1] | Q |

[Add Tags & Comments](#) Comments: 0 [Associated Tags: No Tags](#)

DISTANCE DISTRIBUTION:

40
30

Within-BIN NN-Papilio machaon

BIN COMPLIANT WITH METADATA REQUIREMENTS

Specimen Images:

LEATG387-14 {Papilio machaon}

License: CreativeCommons - Attribution Non-Commercial Share-Alike

License Holder: Peter Huemer, Tiroler Landesmuseum Ferdinandeum

[Add Tags & Comments](#) Comments: 0 [Associated Tags: No Tags](#)

BIN: BOLD:AAA5810

Papilio machaon

3. The metabarcoding approach

METABARCODING



DNA-BARCODING
(2003)

THE ROYAL SOCIETY
Received 29 July 2002
Accepted 30 September 2002
Published online 8 January 2003

Biological identifications through DNA barcodes

Paul D. N. Hebert¹, Alina Cywinski¹, Shelley L. Ball¹ and Jeremy R. deWaard²

¹Department of Zoology, University of Guelph, Guelph, Ontario N1G 2W1, Canada

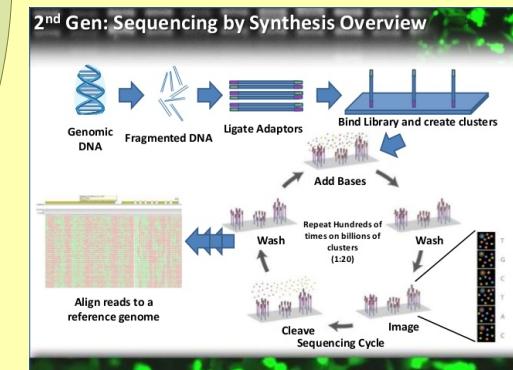
Although much biological research depends upon species diagnoses, taxonomic expertise is collapsing. We are convinced that the sole prospect for a sustainable identification capability lies in the construction of systems that employ DNA sequences as taxon 'barcodes'. We establish that the mitochondrial gene cytochrome c oxidase I (COI) can serve as the core of a global biodiversity system for animals. First, we demonstrate that COI profiles, derived from the low-density sampling of higher taxonomic categories, ordinarily assign newly analysed taxa to the appropriate phylum or order. Second, we demonstrate that species-level assignments can be obtained by creating comprehensive COI profiles. A model COI profile, based upon the analysis of a single individual from each of 200 closely allied species of lepidopterans, was 100% successful in correctly identifying subsequent specimens. When fully developed, a COI identification system will provide a reliable, cost-effective and accessible solution to the current problem of species identification. Its assembly will also generate important new insights into the diversification of life and the rules of molecular evolution.

Keywords: molecular taxonomy; mitochondrial DNA; animals; insects; sequence diversity; evolution

METABARCODING

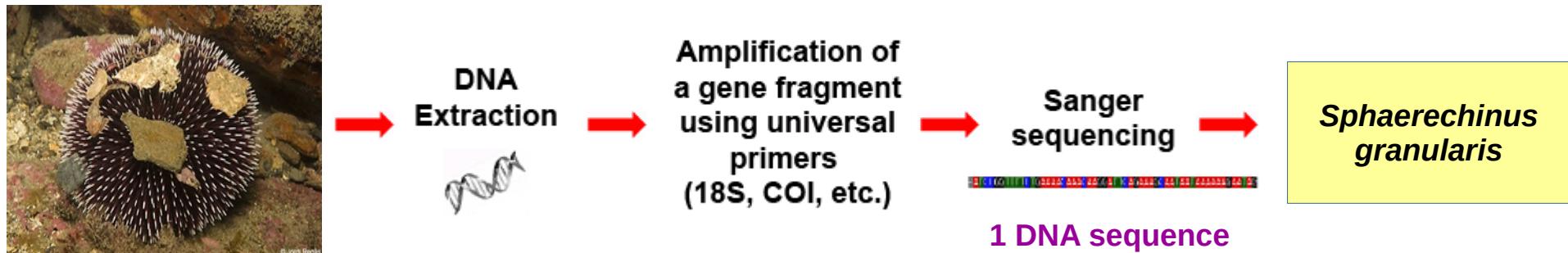


**HIGH-THROUGHPUT
SEQUENCING**
(2002-2010)

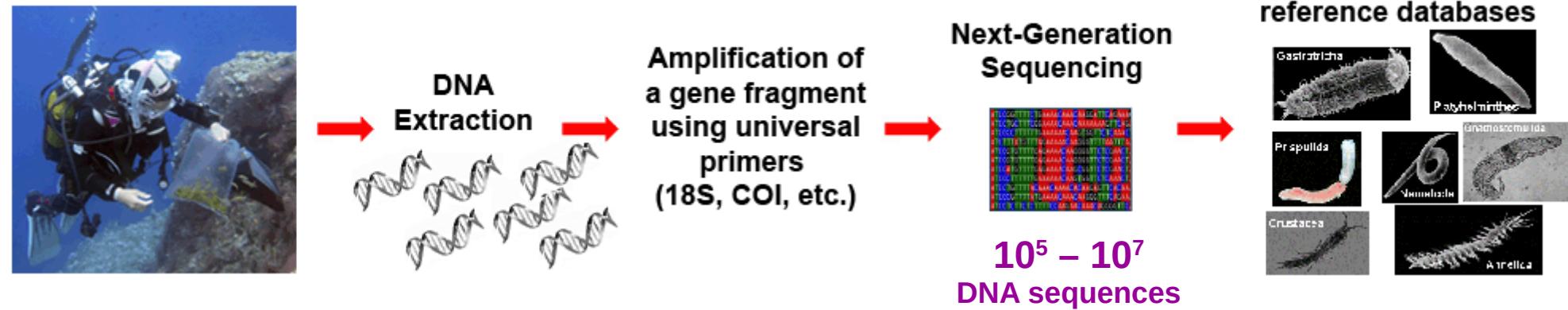


DNA BARCODING vs METABARCODING

DNA BARCODING



DNA METABARCODING



Detects not only living individuals, but also symbionts and parasites, gut contents, body fragments, dead remnants and extracellular DNA

HIGH THROUGHPUT SEQUENCING SHORT FRAGMENTS



Roche GS-FLX+ 454

- Up to 1 M seqs
- 1000 bp max length
- 700 Mb per run



Illumina MiSeq

- Up to 25 M seqs
- 2x 300 bp max length
- 15 Gb per run



Illumina HiSeq 2500

- Up to 600 M seqs
- 2x 250 bp max length
- 300 Gb per run

HIGH THROUGHPUT SEQUENCING SHORT FRAGMENTS



Illumina NextSeq 550

- Up to 400 M seqs
- 2 x 150 bp max length
- 120 Gb per run



Illumina NovaSeq 6000

- Up to 20,000 M seqs
- 2 x 150 bp max length
(2 x 250 bp SP kit)
- 3000 Gb per run

HIGH THROUGHPUT SEQUENCING

LONGER FRAGMENTS



Oxford NT MinION

- Up to 1 M seqs
- Variable. Up to 900 kb length!
- 5 Gb per run

5-15% Error rates!
(for single reads)

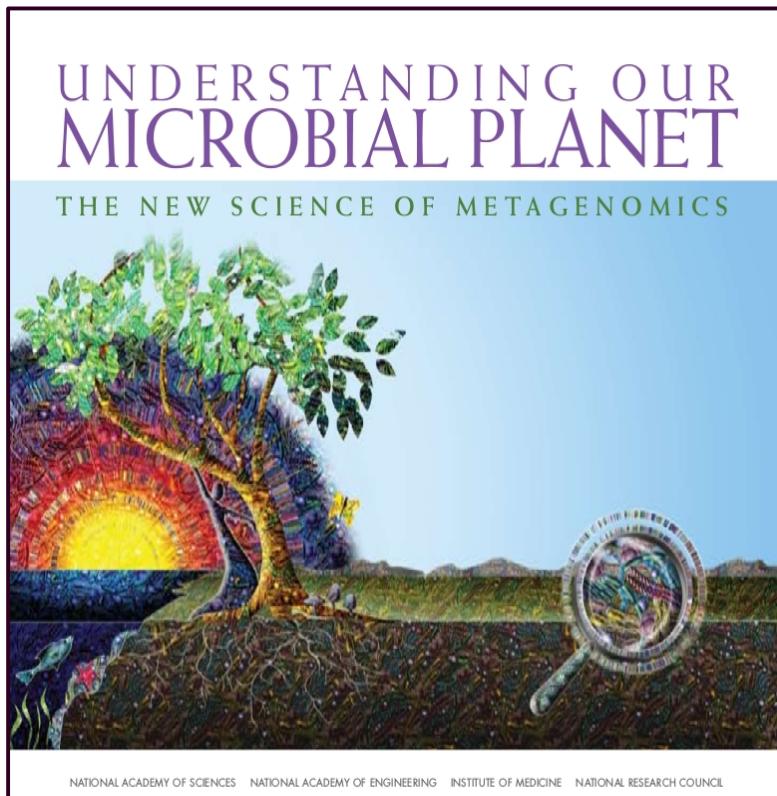


PacBio Sequel

- Up to 500 k seqs
- Up to 35 kb max length
- 10 Gb per run

10-15% Error rates!
(for single reads)

PROKARYOTE 16S METABARCODING

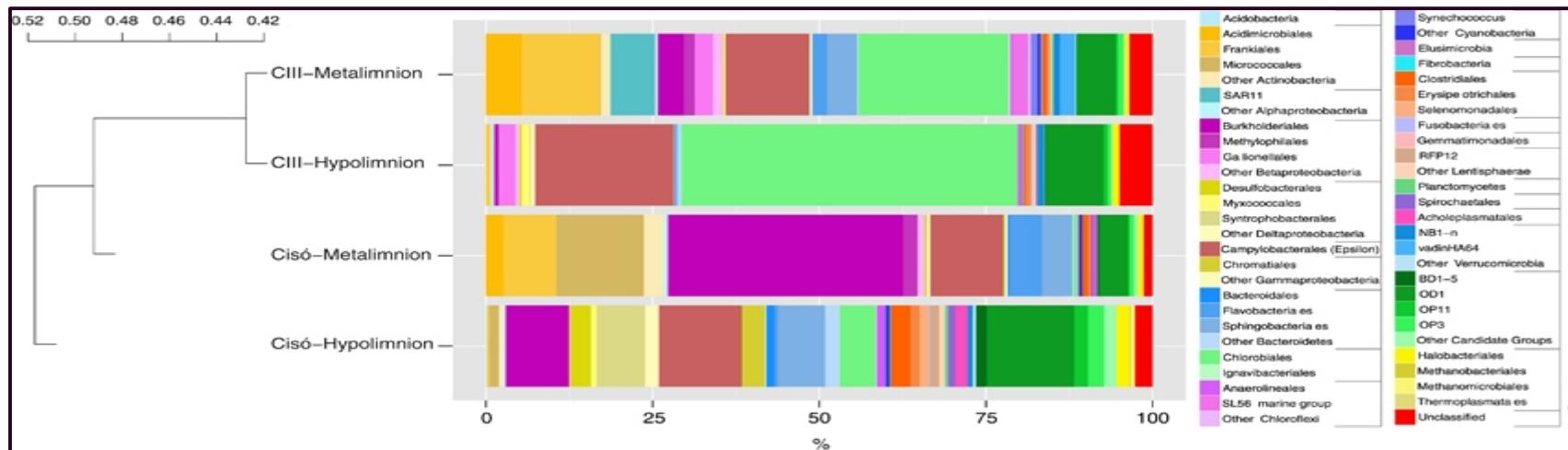


National Academy of Sciences 2007

https://nap.nationalacademies.org/resource/11902/metagenomics_final.pdf

- Human microbiomes
 - Marine microbes
 - Soil prokaryote communities
 - Bioremediation
 - Water treatment ...

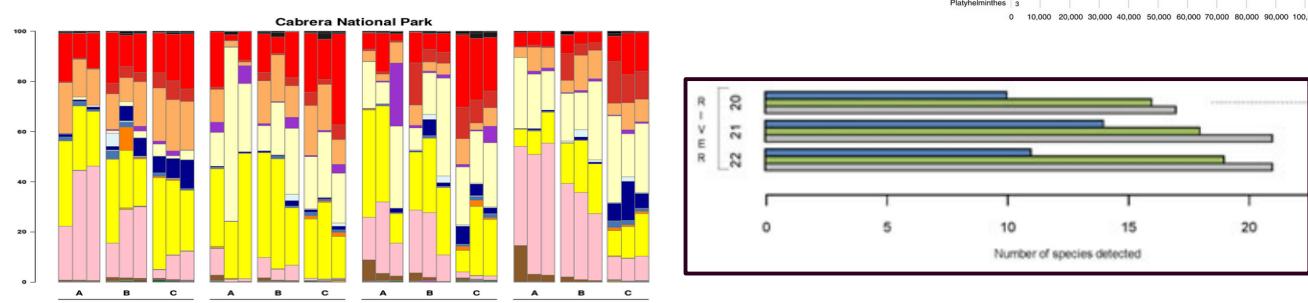
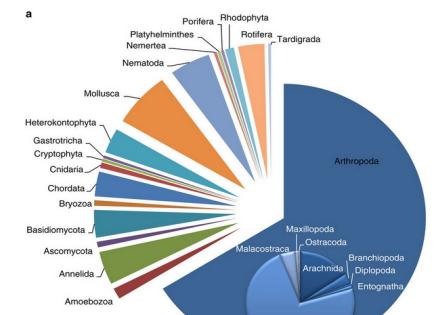
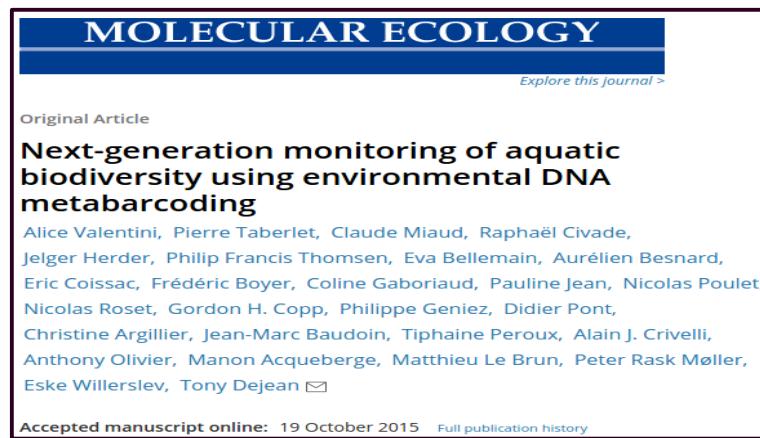
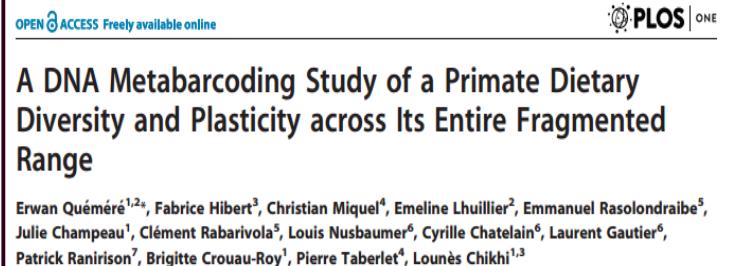
**Well-established pipelines and markers (16S)
Commercial applications (several companies)**



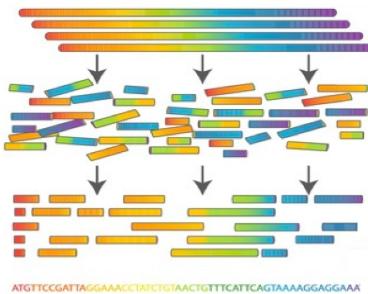
EUKARYOTIC METABARCODING

DEVELOPMENT WORK IN PROGRESS:

- Different sampling techniques, preservation and DNA extraction procedures
- Some confusion between environmental DNA / community DNA / extraorganismal DNA
- Need to decide which markers yield the best results
- Need to find good metabarcoding primers: universal and specific
- Need to standardize analysis pipelines
- Need to improve statistical data treatment

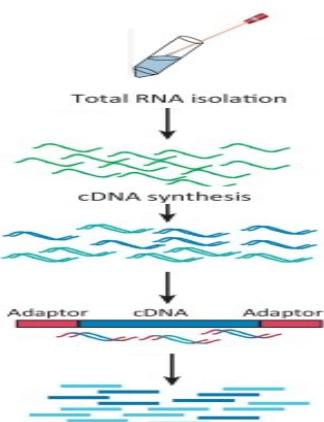


VARIETY OF METAGENOMICS TECHNIQUES



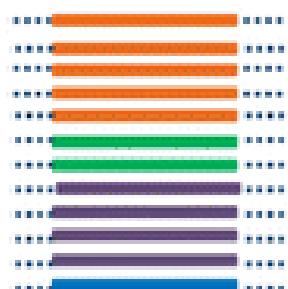
TOTAL DNA SEQUENCING (SHOTGUN SEQUENCING):

- METAGENOMICS (mostly Prokaryotic, currently)
- MITOGENOMICS (MITOCHONDRIAL DNA METAGENOMICS)
- TARGET-ENRICHMENT TECHNIQUES (still experimental)



TOTAL RNA SEQUENCING (cDNA SHOTGUN):

- METATRANSCRIPTOMICS (mostly Prokaryotic, currently)



AMPLICON SEQUENCING (METABARCODING):

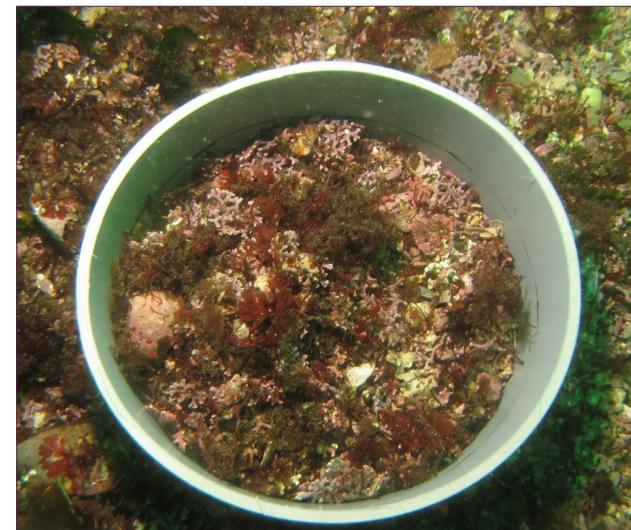
- DNA METABARCODING
- RNA METABARCODING (RIBOSOMAL MARKERS)

TYPES OF GENETIC SAMPLES

ENVIRONMENTAL DNA

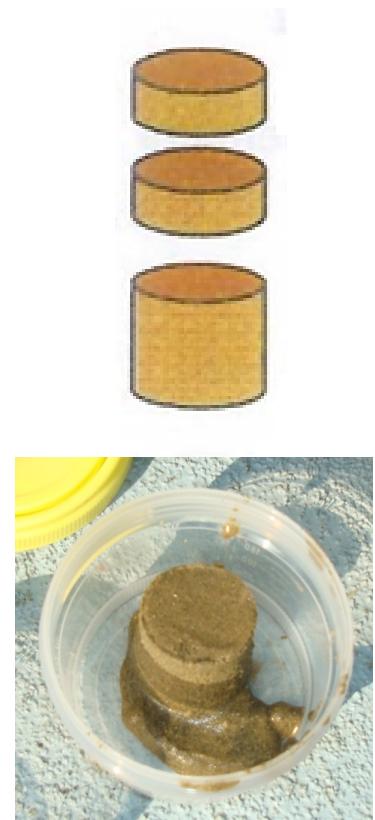


COMMUNITY DNA



THIS DISTINCTION MAY BE BLURRY

MARINE SEDIMENTS



PLANKTON COMMUNITIES

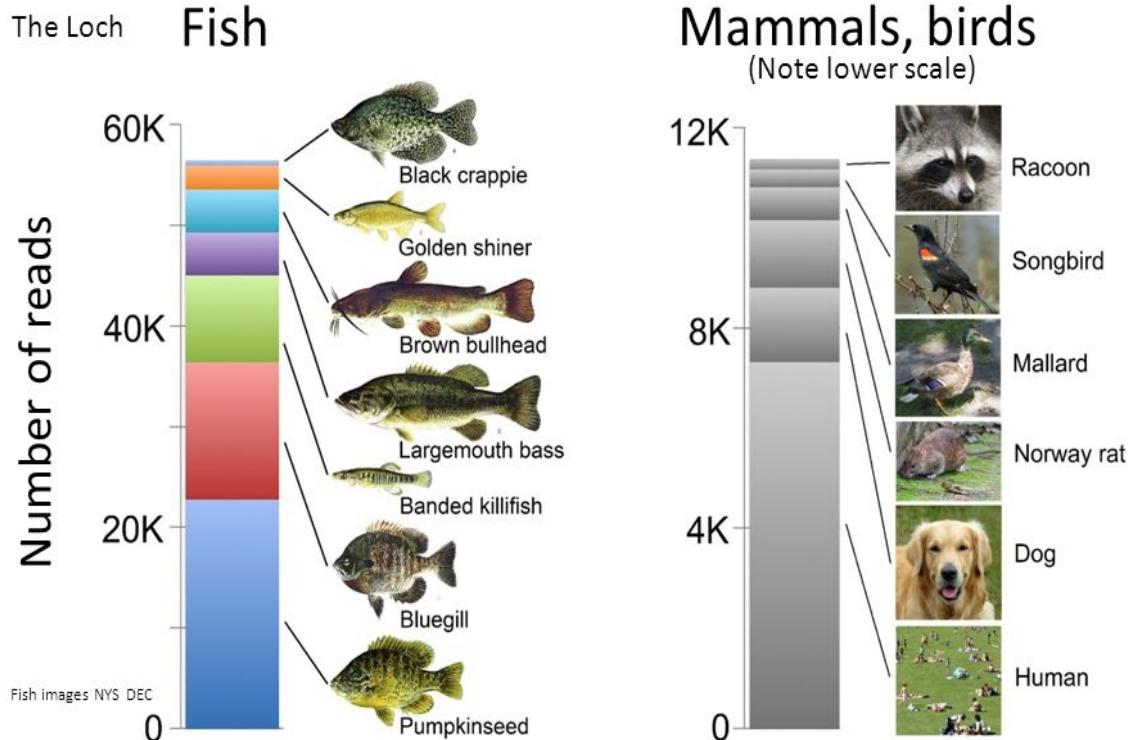


AND DEPENDS ON THE TARGET TAXA (SPECIFICITY OF PRIMERS USED)

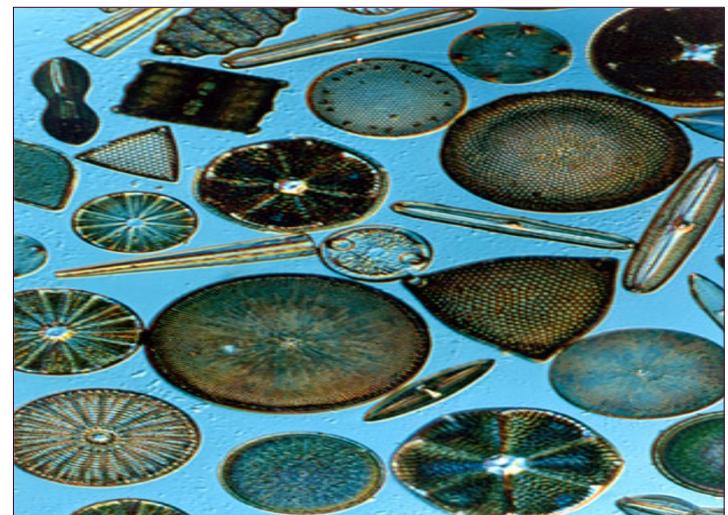
EXTRA-ORGANISMAL DNA

A dozen species in $\frac{1}{4}$ cup of water

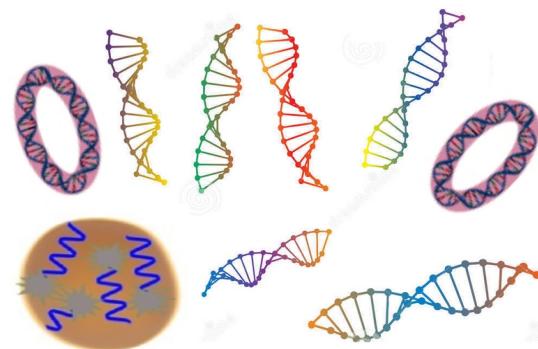
- Analyze 6.33 ng DNA (roughly 60 mL H₂O); 7×10^4 reads
- Detected 7 species of fish, also mammals, birds



COMMUNITY DNA



EVERY eDNA SAMPLE IS ACTUALLY A COMBINATION OF TWO FRACTIONS:



eDNA = extra-organismal DNA + community DNA

- total DNA
- envDNA

- "free" DNA
- "dissolved" DNA
- "true" eDNA
- extracellular DNA
- "dead" DNA
- "trace" DNA

- "living" DNA
- organismal eDNA

Trade-offs between reducing complex terminology and producing accurate interpretations from environmental DNA:
Comment on "Environmental DNA: What's behind the term?"
by Pawlowski et al., (2020)

<https://onlinelibrary.wiley.com/doi/epdf/10.1111/mec.15942>

WITH IMPLICATIONS FOR THE QUANTITATIVE VALUE

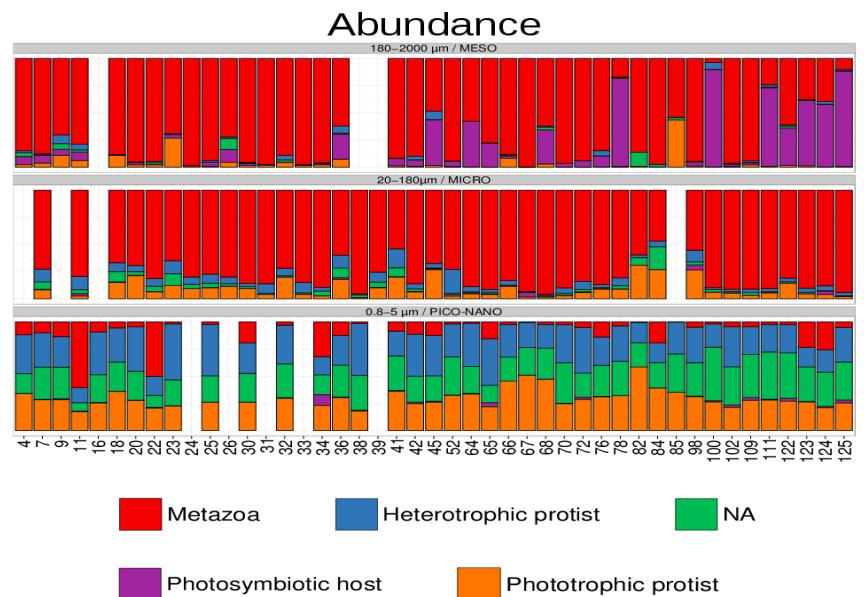
EXTRA-ORGANISMAL DNA



COMMUNITY DNA



PRESENCE / ABSENCE DATA



QUANTITATIVE DATA

SAMPLING STRATEGY: A QUESTION OF TARGET CONCENTRATION



THE NEEDLE?

or

THE HAY?



SUMMARY:

COMMUNITY DNA METABARCODING

- High target DNA concentration
- Deterministic PCR
- Very high repeatability
- No need to use PCR replicates
- A few ecological replicates
- Quantitative value of N. of reads
- Relative abundance of reads reflects relative biomass of MOTUs
- No contamination issues
- Negative controls will be OK

EXTRA-ORGANISMAL DNA METABARCODING

- Very low target DNA concentration
- Stochastic PCR
- Very low repeatability
- Need for many PCR replicates
- As many as possible ecol. replicates
- Stochasticity in N. of reads
- Use percentage of positive PCRs as a measure of MOTU abundance
- Very sensitive to contaminations
- Problematic MOTUs in negatives

COMMUNITY DNA METABARCODING

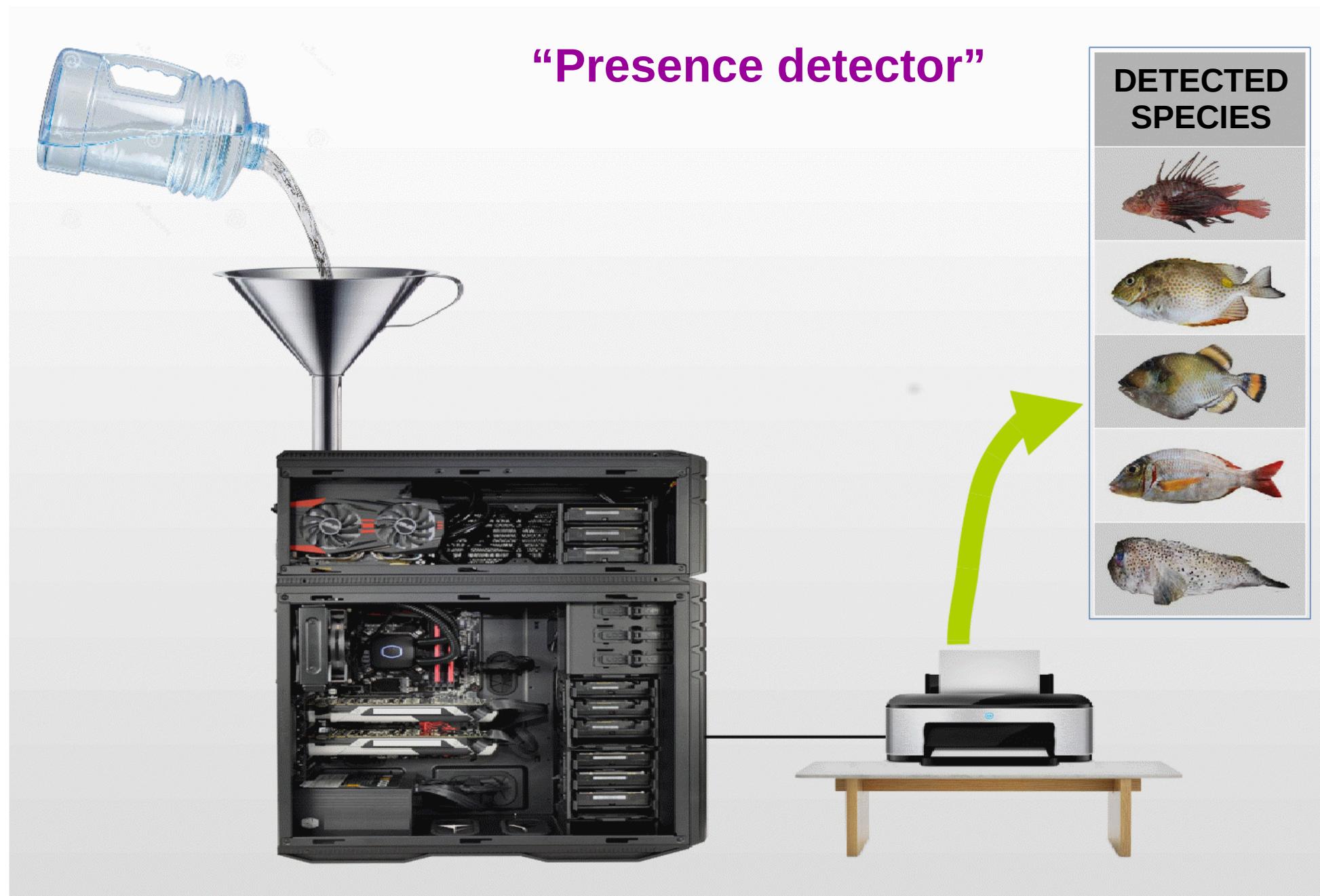


“Automated Inventory Machine”

| Species registered | Natural substrate fauna | | | | | | | | | | | | Nylon mesh-fauna fauna | | | | | | | | | | | | | | |
|---------------------------------------|-------------------------|------------------------------|------------------|-------------|--------|----------|---------------------|--------------------|----------------------|----------------|------------|------------|------------------------|--|--------------------|--|--------------------|--|----------------------|--|--------------------|--|--------------------|---------|----------------|------|------|
| | Type of environment | | | | | | Type of environment | | | | | | Primary Forest (f) | | | | | | Secondary Forest (s) | | | | | | Polyvalent (p) | | |
| | Primary forest (f) | Disturbed primary forest (f) | Cultivations (c) | Savanna (s) | "Iapo" | "Várzea" | Trees (t) | Primary Forest (f) | Secondary Forest (s) | Polyvalent (p) | Forest (f) | Forest (s) | (p) | Number of reference (see Table 1 for more details) | Region of sampling | Number of reference (see Table 1 for more details) | Region of sampling | Number of reference (see Table 1 for more details) | Region of sampling | Number of reference (see Table 1 for more details) | Region of sampling | Number of reference (see Table 1 for more details) | Region of sampling | | | | |
| f1 | f2 | f3 | f4 | f5 | f6 | f7 | f8 | c9 | s10 | s11 | s12 | s13 | s14 | s15 | s16 | s17 | s18 | s19 | s20 | s21 | s22 | s23 | s24 | s25 | s26 | | |
| Peru | RR | RO | PA | AM | AM | AM | AM | RO | AM | BR | RR | RR | PA | PA | AM | AM | AM | AM | AM | AM | AM | AM | AM | AM | AM | | |
| LOWER ORBITIDA | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Level A | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Phalangomorpha | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1.5 | 2.3 | 3.4 | 2.8 | 0 | 0 | 0 | 0 | 0 | 1.2 | 0 | 1.4 | 4.5 | |
| Eupnoptera | 1.3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3.6 | 0 | 0 | 0 | 0 | 8.9 | 11.6 | 0 | 8.3 | 0 | 0 | 0 | 0 | 0 | 3.5 | 0 | 4.2 | 3.0 | |
| Total | 1.3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 7.1 | 0 | 0 | 0 | 0 | 10.5 | 14.0 | 3.4 | 11.1 | 0 | 0.8 | 0 | 0 | 0 | 1.7 | 0 | 0 | 5.8 | 7.6 |
| Level B | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Hypochthoniidae groups | 5.2 | 2.7 | 13.2 | 5.9 | 7.0 | 7.7 | 7.4 | 8.3 | 10.7 | 6.7 | 7.7 | 4.0 | 5.9 | 11.6 | 6.9 | 5.6 | 7.7 | 1.7 | 0 | 0 | 2.9 | 1.4 | 3.5 | 3.6 | 4.3 | 4.5 | |
| Total | 5.2 | 2.7 | 13.2 | 5.9 | 7.0 | 7.7 | 7.4 | 8.3 | 10.7 | 6.7 | 7.7 | 4.0 | 5.9 | 11.6 | 6.9 | 5.6 | 7.7 | 1.7 | 0 | 0 | 2.9 | 1.4 | 3.5 | 3.6 | 4.3 | 4.5 | |
| Level C | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Misnomata | 9.0 | 29.7 | 7.8 | 7.8 | 12.3 | 5.3 | 1.9 | 12.5 | 1.8 | 13.3 | 11.5 | 8.0 | 5.9 | 23 | 17.2 | 2.8 | 0.0 | 5.8 | 4.5 | 0 | 7.4 | 5.9 | 8.1 | 5.5 | 11.6 | 10.6 | |
| Nothroidea s.l. | 5.8 | 10.8 | 10.5 | 9.8 | 8.8 | 12.8 | 11.1 | 4.2 | 5.4 | 6.7 | 7.7 | 12.0 | 4.5 | 16.3 | 3.4 | 13.9 | 15.4 | 3.3 | 4.5 | 0 | 5.9 | 8.5 | 10.5 | 7.3 | 8.7 | 9.1 | |
| Total | 14.8 | 40.5 | 18.4 | 17.6 | 21.4 | 17.9 | 13.0 | 16.7 | 7.1 | 20.0 | 19.2 | 1.5 | 18.6 | 20.7 | 16.7 | 15.4 | 9.2 | 9.1 | 0 | 13.2 | 18.3 | 18.6 | 12.7 | 20.3 | 19.7 | 21.3 | |
| Total Lower Orbitida | 21.3 | 43.2 | 31.6 | 21.5 | 28.1 | 25.8 | 20.4 | 25.0 | 25.0 | 26.9 | 24.0 | 26.9 | 44.2 | 31.0 | 33.3 | 28.1 | 11.7 | 9.1 | 0 | 16.2 | 19.7 | 26.7 | 16.4 | 30.4 | 31.8 | | |
| HIGHER ORBITIDA | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Acaroidea organization level A | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Ancient Opilio and Eupnoptera | 6.5 | 0 | 2.6 | 5.9 | 5.3 | 3.8 | 3.7 | 8.3 | 5.4 | 0 | 0 | 4.0 | 3.0 | 2.3 | 3.4 | 0 | 0 | 4.2 | 4.5 | 25.0 | 4.4 | 5.6 | 7.0 | 7.3 | 4.3 | 6.1 | |
| Cytheromorpha and related groups | 0 | 2.7 | 0 | 0 | 0 | 1.3 | 0 | 0 | 1.8 | 0 | 0 | 0 | 0 | 3.0 | 0 | 0 | 2.8 | 0 | 1.7 | 4.5 | 25.0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Carabodesidae | 5.2 | 0 | 2.6 | 3.9 | 8.8 | 11.5 | 18.5 | 4.2 | 12.5 | 0 | 0 | 4.0 | 0 | 0 | 0 | 2.8 | 0 | 5.0 | 9.1 | 0 | 7.4 | 5.6 | 11.6 | 7.3 | 11.6 | 6.1 | |
| Total | 11.6 | 2.7 | 5.3 | 9.8 | 14.0 | 16.6 | 22.2 | 12.5 | 19.6 | 0 | 0 | 0 | 0 | 8.0 | 6.0 | 2.3 | 3.4 | 5.6 | 0 | 10.8 | 18.2 | 50.0 | 11.8 | 13.3 | 18.6 | 14.5 | 12.3 |
| Gymnophoroidae | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Archopagomerinae | 3.9 | 2.7 | 0 | 2.0 | 1.8 | 6.4 | 5.6 | 0 | 5.4 | 6.7 | 7.7 | 0 | 4.5 | 4.7 | 3.4 | 11.1 | 0 | 3.3 | 0 | 0 | 4.4 | 2.8 | 2.3 | 1.8 | 1.4 | 0 | |
| Lioscapididae and related genera | 3.9 | 0 | 0 | 3.8 | 0 | 2.6 | 7.4 | 0 | 7.1 | 0 | 0 | 0 | 0 | 2.3 | 0 | 2.8 | 0 | 0 | 4.5 | 0 | 2.9 | 1.4 | 2.3 | 1.8 | 4.3 | 1.5 | |
| Total | 6.5 | 2.7 | 0 | 5.9 | 3.5 | 9.0 | 13.0 | 0 | 14.3 | 6.7 | 11.5 | 0 | 6.0 | 11.6 | 3.4 | 16.7 | 0 | 5.0 | 4.5 | 0 | 7.4 | 4.2 | 4.7 | 3.6 | 5.8 | 1.5 | |
| Advanced level C | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Brachynellidae | 7.1 | 5.4 | 13.2 | 7.8 | 5.3 | 10.2 | 14.8 | 8.3 | 12.5 | 13.3 | 7.7 | 8.0 | 3.0 | 9.2 | 10.3 | 11.0 | 7.7 | 4.2 | 9.1 | 0 | 8.4 | 11.3 | 7.8 | 5.8 | 9.1 | | |
| Basic almost apertic paleopagomerinae | 7.1 | 0 | 2.6 | 0 | 5.3 | 10.3 | 3.7 | 12.5 | 3.6 | 0 | 3.8 | 8.0 | 9.0 | 7.0 | 0 | 5.6 | 0 | 5.0 | 13.6 | 0 | 7.4 | 11.3 | 7.0 | 9.1 | 11.6 | 6.1 | |
| Basic poromidae | 2.6 | 2.7 | 2.6 | 2.0 | 1.8 | 5.1 | 3.7 | 0 | 8.9 | 0 | 0 | 0 | 0 | 2.3 | 3.4 | 8.3 | 0 | 4.2 | 4.5 | 0 | 1.5 | 2.8 | 1.2 | 1.8 | 0 | 1.5 | |
| Periperic poromidae paleopagomerinae | 5.8 | 2.7 | 5.3 | 5.9 | 10.5 | 7.7 | 9.3 | 8.3 | 7.1 | 13.3 | 7.7 | 4.0 | 4.5 | 7.0 | 6.9 | 5.6 | 23.1 | 12.5 | 4.5 | 25.0 | 11.8 | 8.5 | 9.3 | 18.10.1 | 9.1 | | |
| Principia | 5.8 | 8.1 | 5.3 | 2.0 | 5.3 | 0 | 0 | 4.2 | 0 | 13.3 | 11.5 | 0 | 4.5 | 0 | 3.4 | 0 | 0 | 5.8 | 0 | 0 | 7.4 | 2.8 | 3.5 | 9.1 | 2.9 | 6.1 | |
| Total | 28.4 | 18.9 | 28.9 | 17.6 | 28.1 | 42.3 | 31.5 | 33.3 | 32.1 | 40.0 | 30.8 | 20.0 | 20.9 | 25.6 | 24.1 | 33.3 | 30.8 | 11.7 | 31.8 | 25.0 | 36.8 | 36.6 | 26.7 | 29.1 | 30.4 | 31.8 | |
| Terminal level D | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| (Derived Higher Orbitida) | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Derived Aphelinomorpha | 20.6 | 24.3 | 18.4 | 25.5 | 22.8 | 6.4 | 13.0 | 29.2 | 5.4 | 26.7 | 26.9 | 32.0 | 19.4 | 7.0 | 27.6 | 8.3 | 46.2 | 28.3 | 31.8 | 25.0 | 16.2 | 12.7 | 14.0 | 16.4 | 10.1 | 12.3 | |
| (Coccoidea) | 11.6 | 8.1 | 15.8 | 17.6 | 3.5 | 0 | 0 | 0 | 3.6 | 0 | 3.8 | 16.0 | 20.9 | 9.3 | 10.3 | 2.8 | 0 | 12.5 | 4.5 | 0 | 11.8 | 15.5 | 9.3 | 20.0 | 7.2 | 10.6 | |
| Total | 32.3 | 32.4 | 34.3 | 43.1 | 26.3 | 6.4 | 13.0 | 29.2 | 8.9 | 26.7 | 30.3 | 48.0 | 40.3 | 16.3 | 37.9 | 11.1 | 46.2 | 40.8 | 36.4 | 25.0 | 27.0 | 28.2 | 33.3 | 36.4 | 17.4 | 22.7 | |
| Total Higher Orbitida | 78.7 | 56.7 | 68.4 | 76.5 | 71.9 | 74.4 | 79.6 | 75.0 | 75.0 | 73.3 | 73.1 | 76.0 | 73.1 | 55.8 | 69.7 | 76.9 | 88.3 | 90.9 | 100.0 | 83.8 | 80.3 | 73.3 | 83.6 | 69.6 | 68.2 | | |

* Relative amount of species within the total number of species.

Extra-organismal DNA METABARCODING

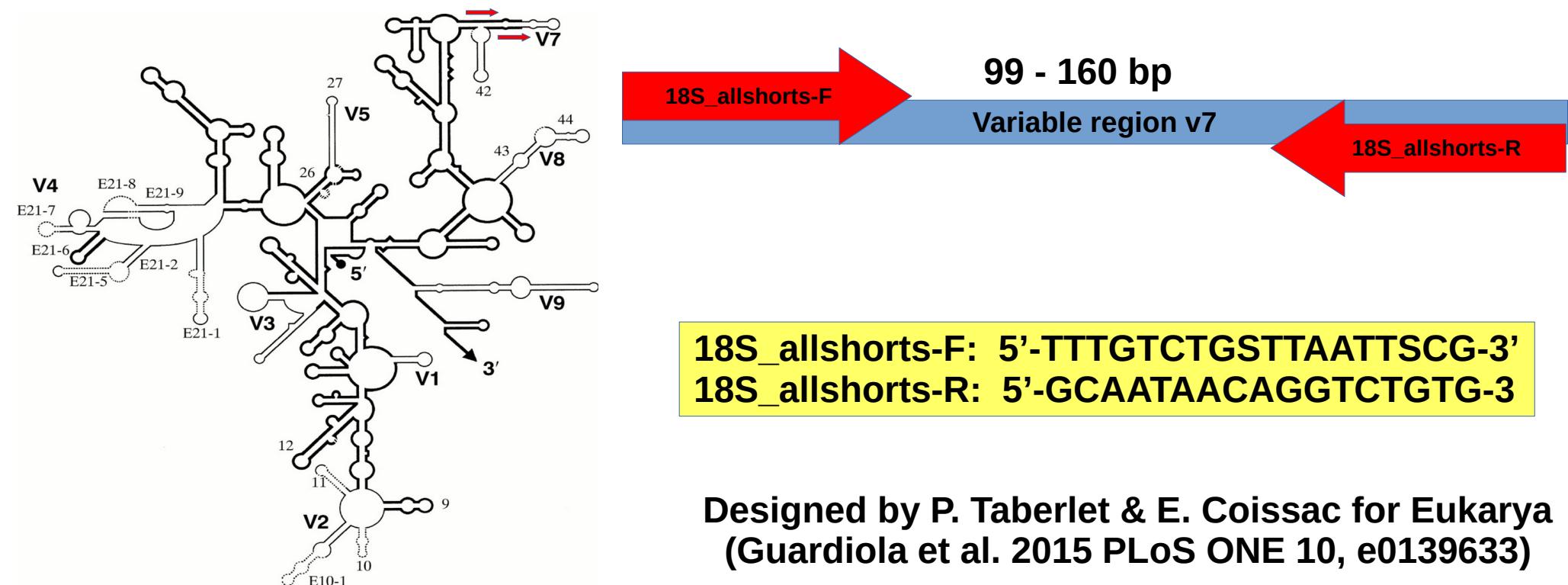


PRIMER DESIGN

METABARCODING MARKER SELECTION:

- A region of enough variability, surrounded by universally conserved regions
- Must be the right size for High-Throughput sequencing (<400 bp in Illumina)
- Ideally a multiple-copy, high-abundance genomic region

EXAMPLE: 18S rRNA v7 region



METABARCODING MARKERS

METABARCODING MARKER SELECTION:

Important issues to consider:

- Variability in length (coding vs non-coding regions)
- Taxon resolution (enough variability to distinguish among taxa?)
- Primer specificity (do the primers amplify only the target taxa?)
- Primer universality (do the primers amplify ALL the target taxa?)
- Intracellular localization (nuclear, mitochondrial, chloroplastidic?)
- Genomic copy number (is it conserved among taxa?)

PRIMER DESIGN AND PCR

METABARCODING MARKER SELECTION

Different metabarcoding markers may be used for different taxonomic groups

The taxonomic target depends on the universality/specifity of the primers

PROKARYOTA: rRNA 16S

METAZOA: nuclear rRNA18S, mitochondrial COI

Vertebrates: mitochondrial rRNA12S

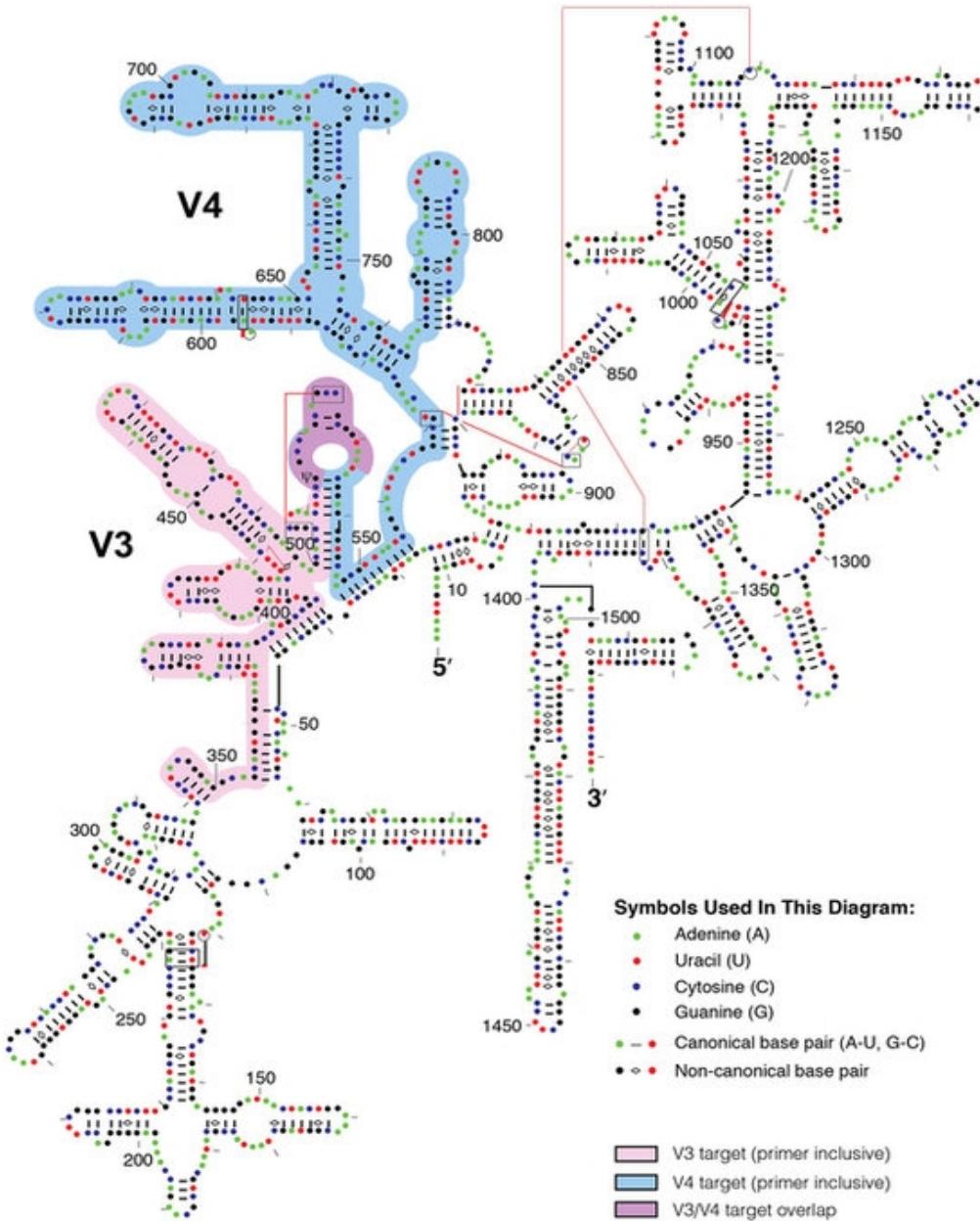
“Invertebrates”: mitochondrial rRNA16S

Arthropoda: mitochondrial COI

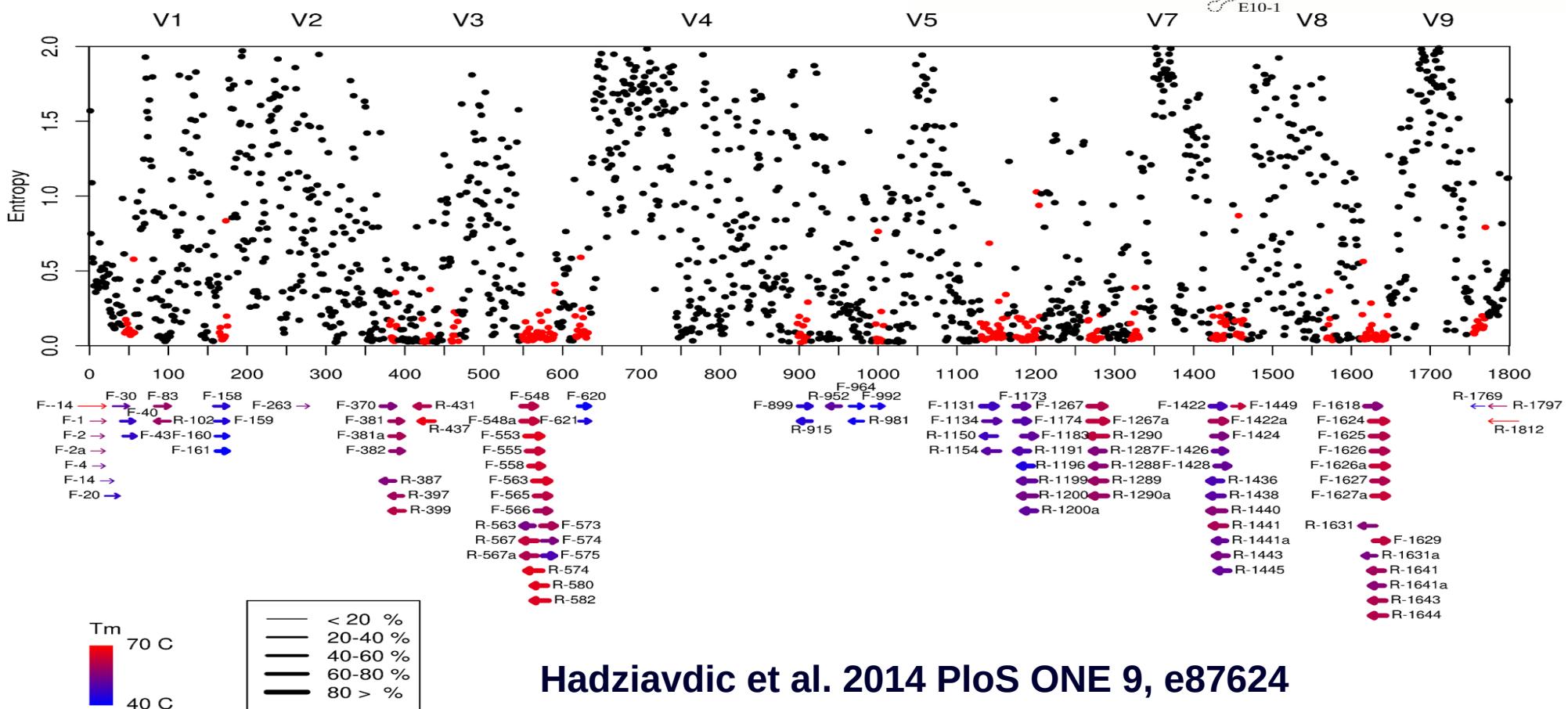
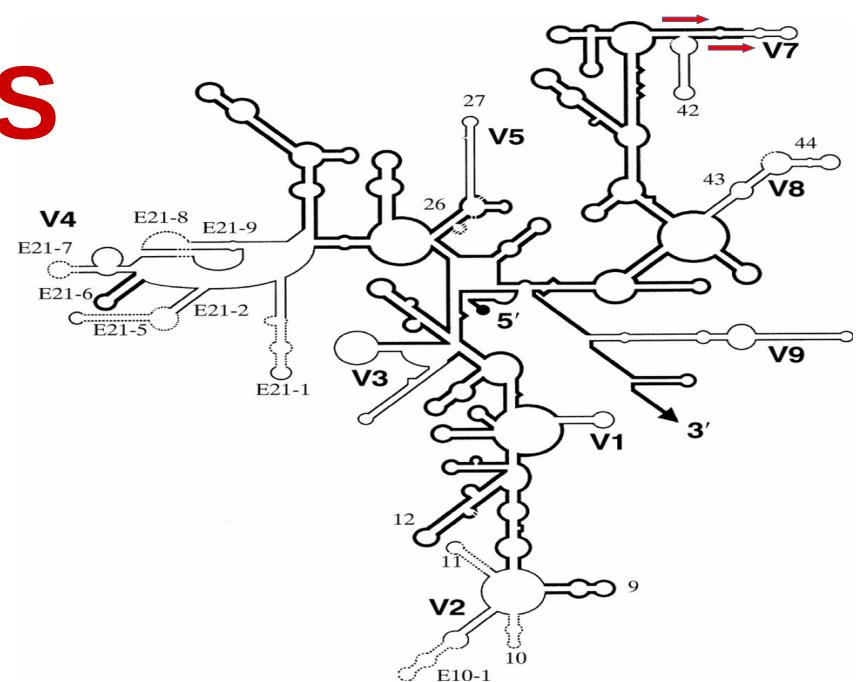
FUNGI: nuclear ITS

PHOTOSYNTHETIC ORGANISMS: chloroplastic trnL, matK, rbcL,
nuclear ITS

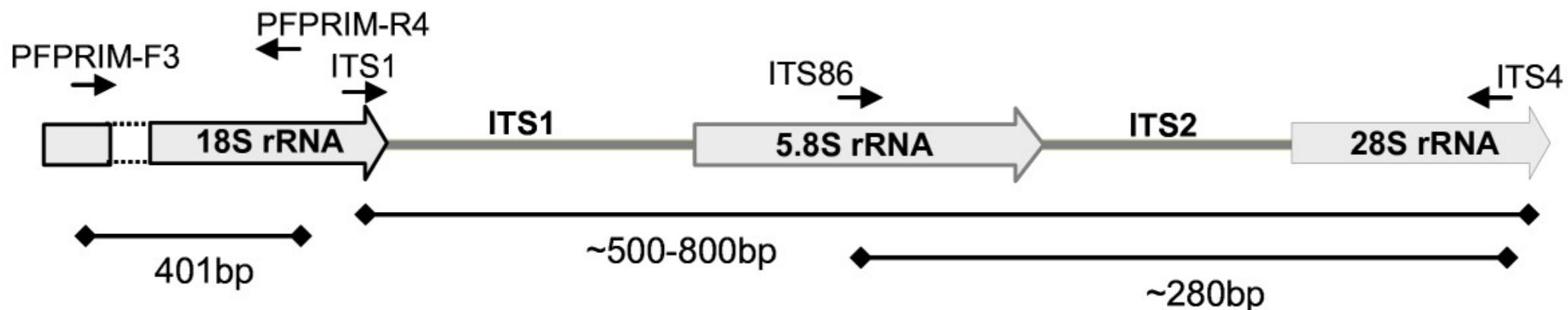
PROKARYOTIC rRNA 16S



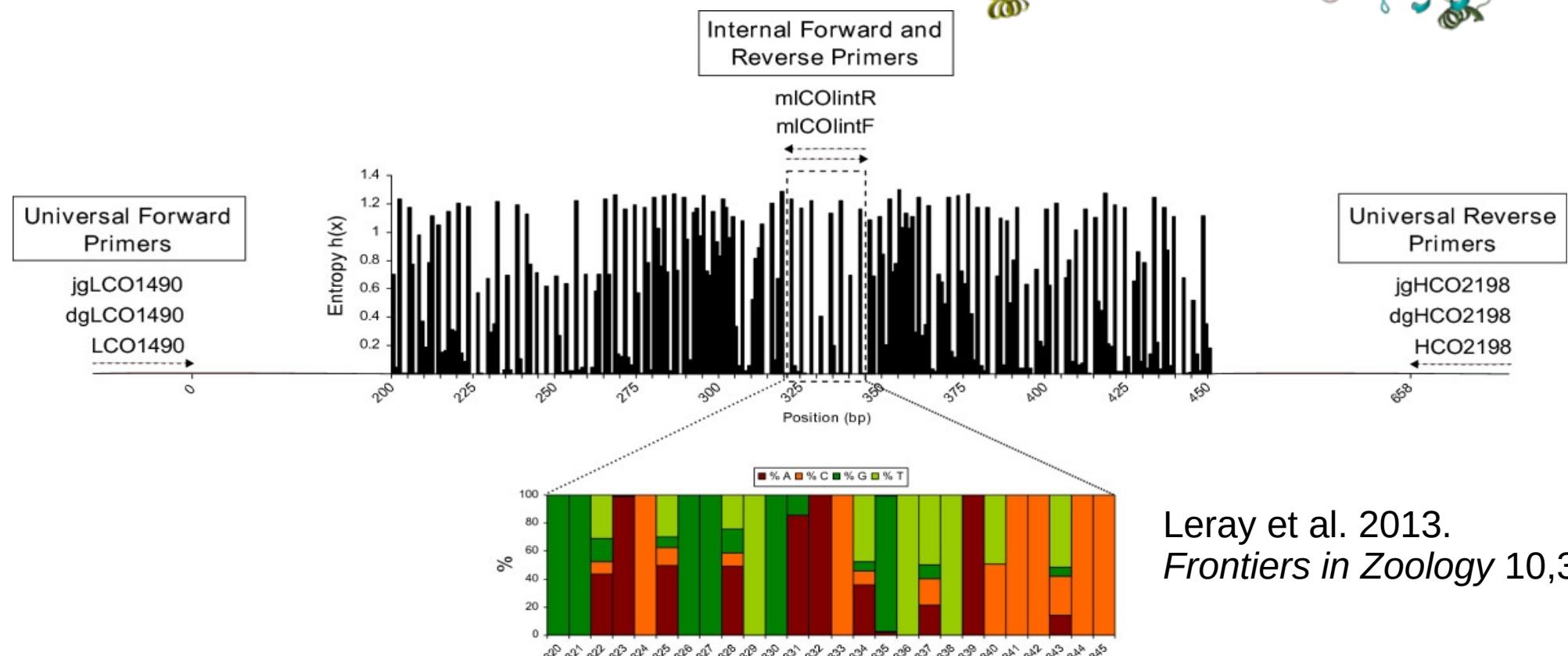
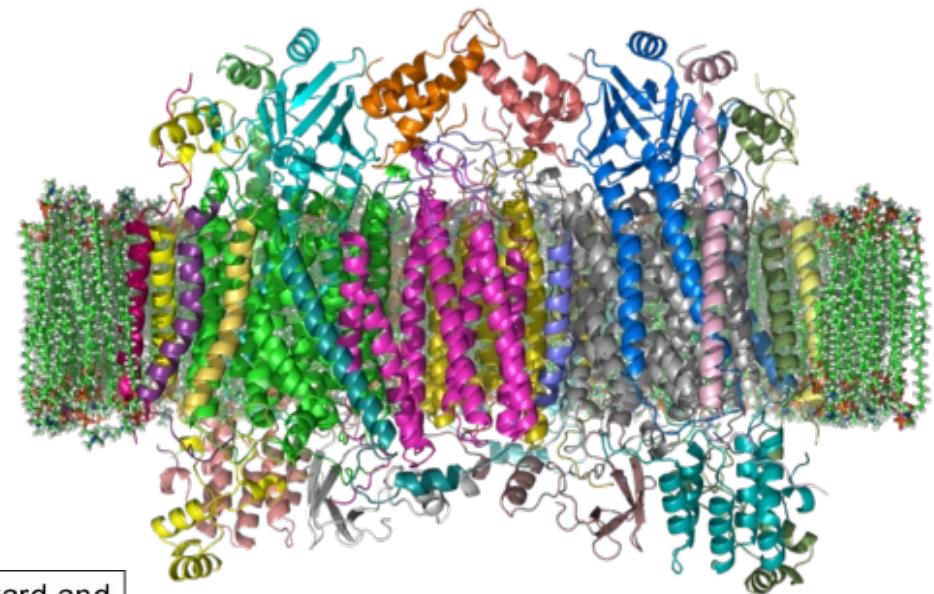
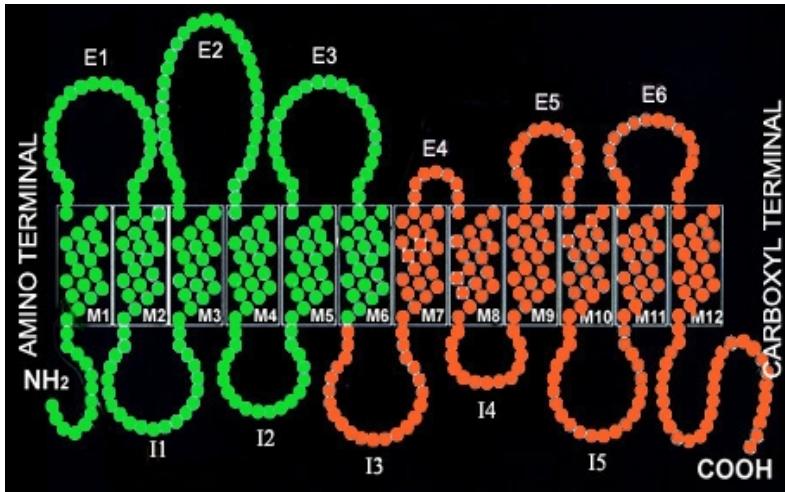
EUKARYOTIC rRNA 18S



INTERNAL TRANSCRIPT SPACERS ITS1 & ITS2



CYTOCHROME c OXIDASE SUBUNIT I



PRIMER DESIGN AND PCR

FREQUENTLY USED METABARCODING MARKERS:

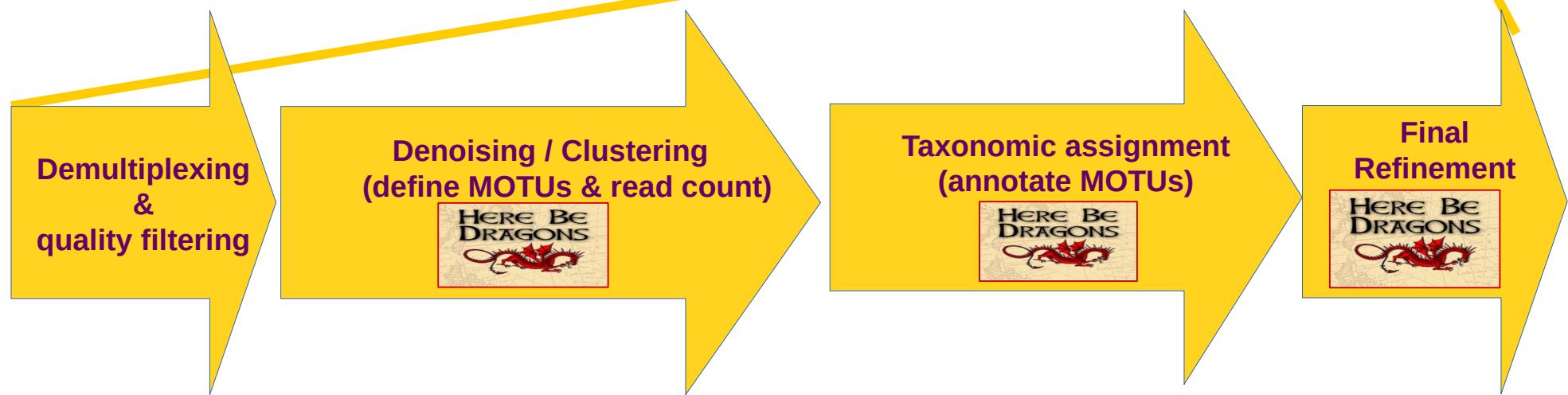
| Target | Gene/ region | Reference | Data bases |
|------------|-----------------|--|--|
| Bacteria | 16S | Sogin <i>et al.</i> (2006) | RDP, Greengenes, SILVA |
| Archaea | 16S | Sogin <i>et al.</i> (2006) | RDP, Greengenes, SILVA |
| Fungi | ITS | Epp <i>et al.</i> (2012), Schoch <i>et al.</i> (2012) | UNITE, GenBank, BOLD (few records) |
| | 18S | Not recommended (Schoch <i>et al.</i> 2012) | SILVA |
| Protists | 18S | Pawlowski <i>et al.</i> (2012) | SILVA |
| | ITS | Pawlowski <i>et al.</i> (2012) | GenBank |
| | CO1 | Pawlowski <i>et al.</i> (2012) | BOLD |
| Meiofauna | CO1 | Hebert <i>et al.</i> (2003) | BOLD |
| | 18S | Deagle <i>et al.</i> (2014) | GenBank |
| Macrofauna | CO1 | Hebert <i>et al.</i> (2003) | BOLD |
| | 16S | Epp <i>et al.</i> (2012), Deagle <i>et al.</i> (2014) | GenBank |
| | 12S | Epp <i>et al.</i> (2012), Deagle <i>et al.</i> (2014) | GenBank |
| | 18S | Deagle <i>et al.</i> (2014) | GenBank |
| Plants | matK | Hollingsworth <i>et al.</i> (2009) | GenBank, BOLD (few records) |
| | + rbcL | | |
| | ITS | China Plant Barcode of Life Group (2011) | GenBank |

4. Metabarcoding workflow

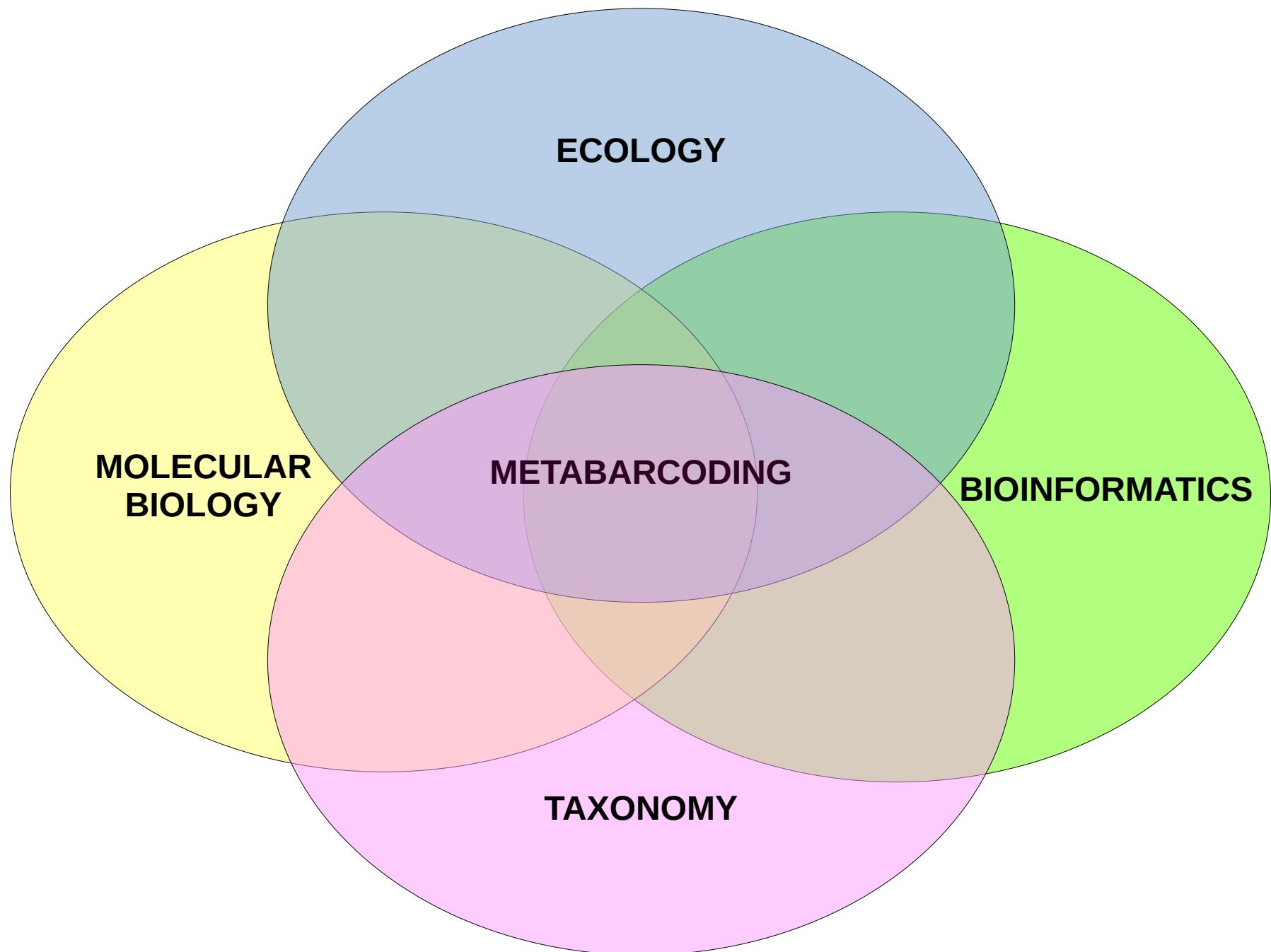
THE METABARCODING WORKFLOW

1. Experimental design, sampling & sample preservation
2. Sample pre-processing
3. DNA extraction
4. DNA amplification (metabarcoding primers)
5. High-throughput sequencing
6. Bioinformatic data analysis
 - 6.1. Demultiplexing and quality control
 - 6.2. Clustering sequences into MOTUs
(Molecular Operational Taxonomic Units)
 - 6.3. Taxonomic assignment
 - 6.4. Refinement of the final dataset
7. Ecological & statistical data treatment

THE METABARCODING WORKFLOW



DISCIPLINES AND SKILLS



DEMULTIPLEXING AND QUALITY FILTERING

Typical MiSeq output: 15 - 25 million raw reads (sequences)

Bioinformatic tools must:

- Remove low-quality sequences
- Assign every sequence to its sample (demultiplexing)
- Group all equal sequences (dereplicating), keeping information of the samples where they occur
- Delete chimeric sequences

Still a huge amount of unique sequences!!!

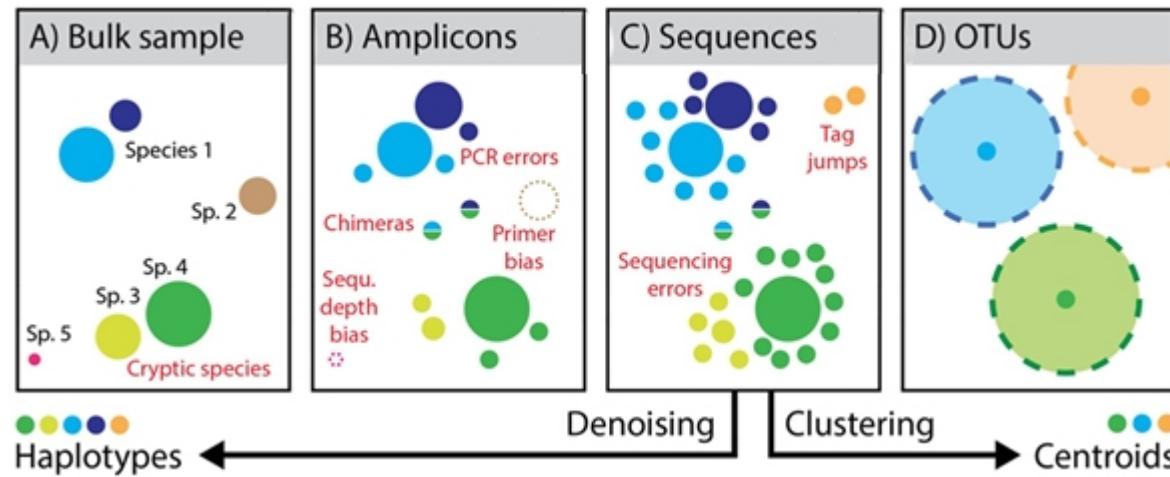
Typically: 1 - 4 million unique sequences

Due to:

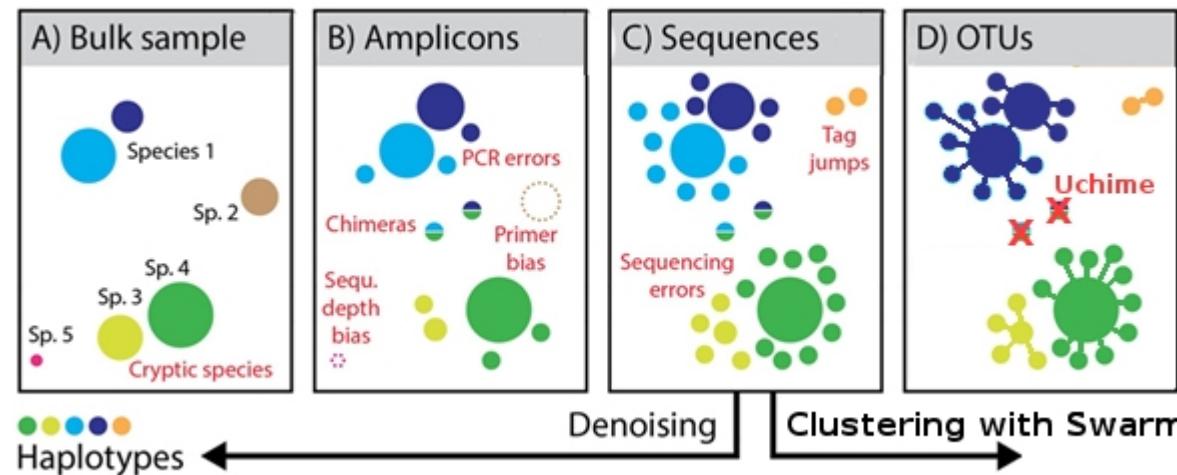
- diversity of species in the samples
- natural intraspecific variability
- random PCR errors
- random sequencing errors

WHAT TO DO WITH ERROR-CONTAINING READS?

CLUSTERING vs DENOISING



Constant vs variable cutoff clustering



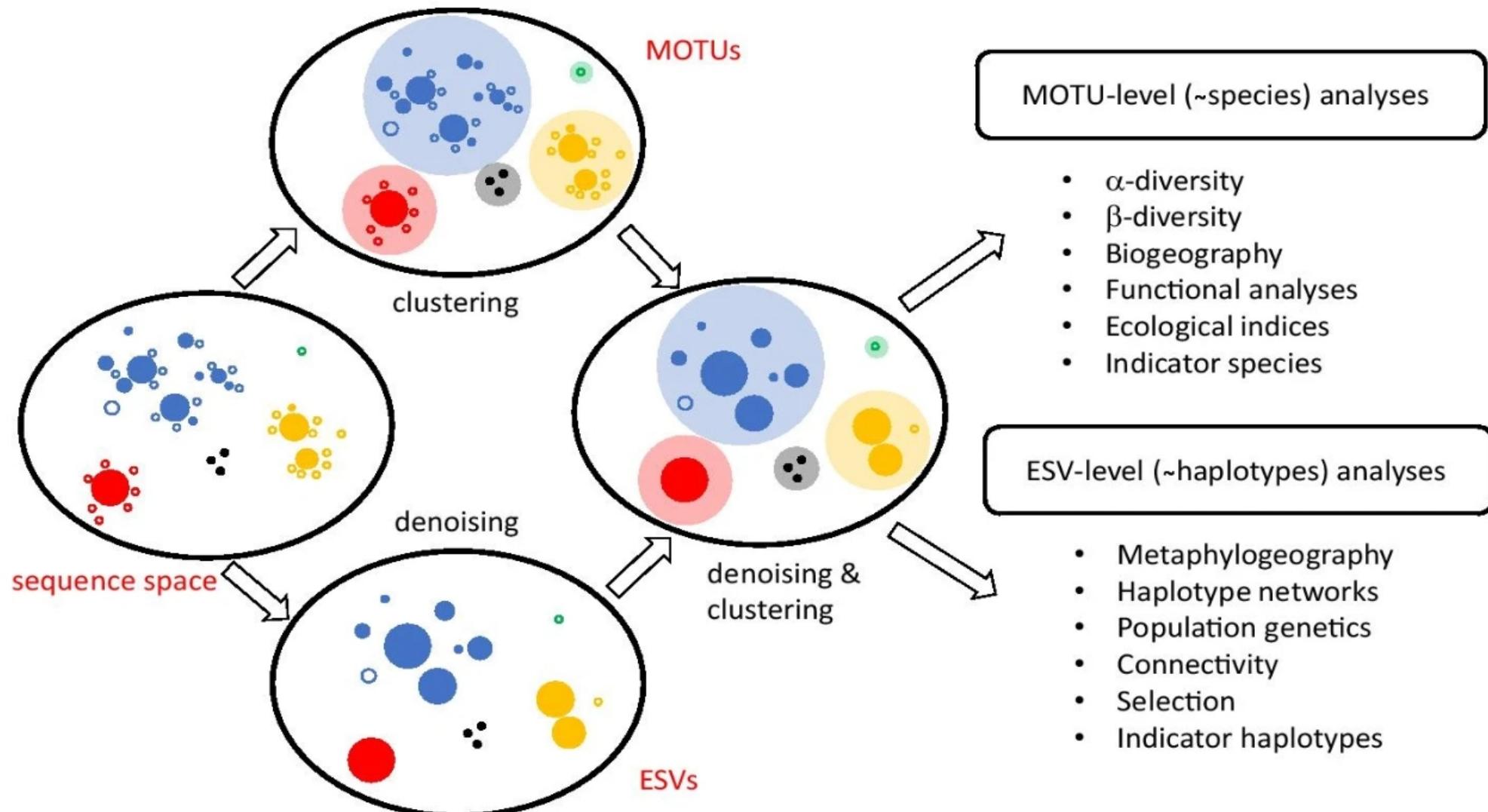
INTEGRATING CLUSTERING & DENOISING

To denoise or to cluster? That is not the question. Optimizing pipelines for COI metabarcoding and metaphylogeography

✉ A. Antich, ✉ C. Palacin, ✉ O.S. Wangensteen, ✉ X. Turon

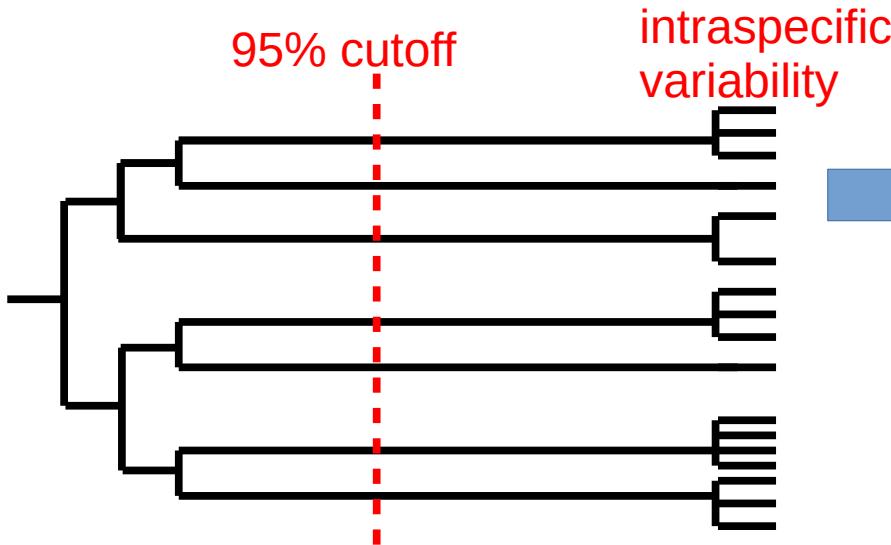
<https://doi.org/10.1101/2021.01.08.425760>

- Denoising AND clustering
(not just denoising OR just clustering)



MOTU CLUSTERING

Theory: Constant ID % cutoff



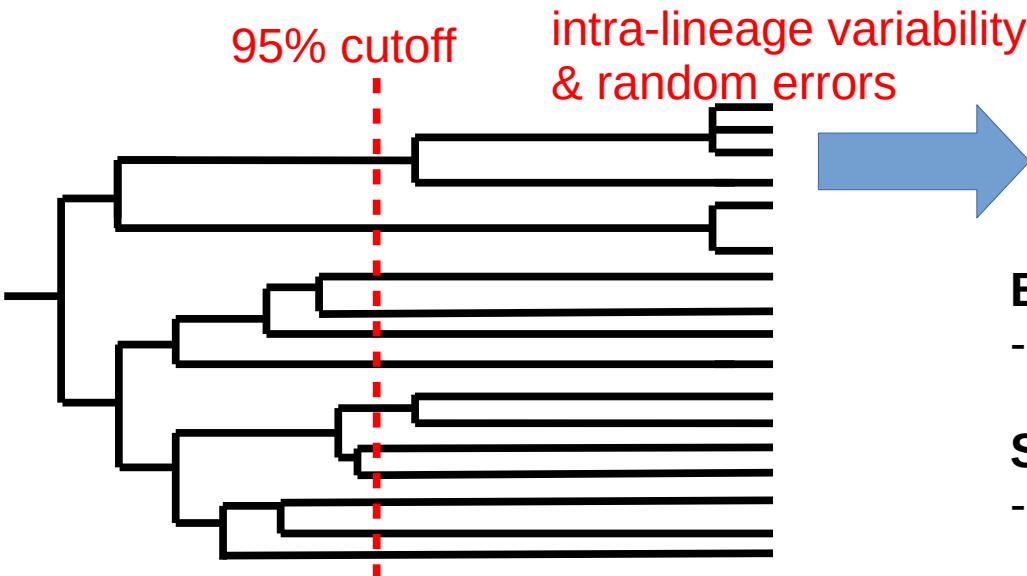
17 different sequences

7 MOTUs (species)

Constant cutoff clustering software:

- UCLUST (Edgar, 2010)
- VSEARCH (Rognes et al 2016)
- SUMACLUST (Mercier et al., 2013)

Reality: Variable cutoff



12 MOTUs (species?)

Bayesian Clustering Algorithms:

- CROP (Hao, Jiang & Chen, 2011)

Step-by-step Aggregation Algorithms:

- SWARM (Mahé et al. 2014)

BIOINFORMATICS PIPELINE: TAXONOMIC ASSIGNMENT OF MOTUS

Ideal: assignment by direct match (ID > 97%)

Marthasterias glacialis voucher LMBRUI7-001 cytochrome oxidase subunit 1 (COI) gene, partial cds; mitochondrial
Sequence ID: [gb|KF369150.1](#) Length: 658 Number of Matches: 1

| Range 1: 346 to 658 GenBank Graphics | | | | | ▼ Next Match | ▲ Previous Match |
|--|--|--------------|-----------|-----------|--------------|------------------|
| Score | Expect | Identities | Gaps | Strand | | |
| 556 bits(616) | 7e-155 | 311/313(99%) | 0/313(0%) | Plus/Plus | | |
| Query 1 | TCTCTCTAGCGGATTAGCCATGCTGGAGGATCAGTGGACCTCGCTATATTCTCTTCA | 60 | | | | |
| Sbjct 346 | TCTCTCTAGCGGATTAGCCATGCTGGAGGATCAGTAGACCTCGCTATATTCTCTTCA | 405 | | | | |
| Query 61 | CTTAGCTGGTGCCTCATCAATCCTTGCCCTCATAAAAATTATCACAAACAATAATTAT | 120 | | | | |
| Sbjct 406 | CTTAGCTGGTGCCTCATCAATCCTTGCCCTCATAAAAATTATCACAAACAATAATTAT | 465 | | | | |
| Query 121 | GCGAACCTCTGGCATGTCCATTGACCGCTTACCCCTATTGTTGATCAGTCTTGTAAAC | 180 | | | | |
| Sbjct 466 | GCGAACCTCTGGCATGTCCATTGACCGCTTACCCCTATTGTTGATCAGTCTTGTAAAC | 525 | | | | |

- BOLD identification
- BLAST against Genbank
- Use local databases (obtained from *in silico* PCR)

What to do with the Unmatched? Phylogenetic Assignment

- Statistical Assignment Package (SAP) Bayesian phylogenetics or constrained-NJ (Munch et al., 2008)
- OBItools - Ecotag Supervised classification algorithm (Boyer et al., 2016)
- PROTAX Probabilistic taxonomic assignment (Somervuo et al. 2016)

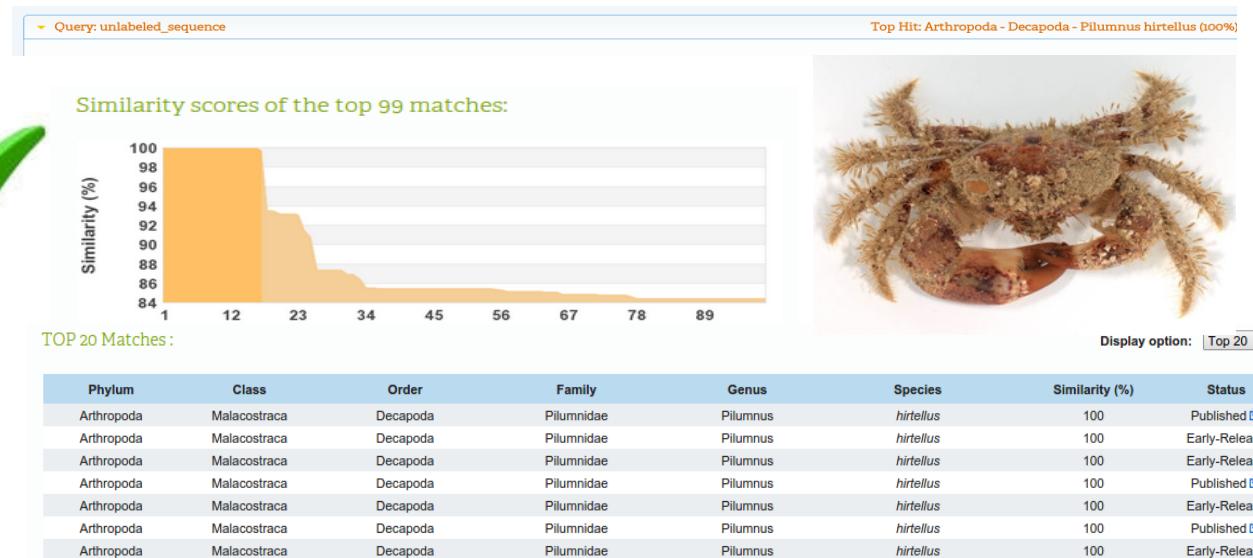
| | | | |
|---------------------------------------|----------------------------|----------------------|---------------------|
| HWI-M00234:272:00000000-AE...: | 1.00 ATTATCAAGAAAATTGCTC | ATGCTGGAAGATCTGTTGAC | CTAGCAAATTTCCTCTAC |
| IN543468_Bohadschia_sp: | 0.82 CTTATCTAGAAAAATTAGCCC | ATGCTGGGGATCAGTAGAC | CTTGCACATTTCCTTACA |
| FJ589208_Actinopyga_lecanora: | 0.82 ACTATCTAGAAAAATTAGCTC | ACGCCAGGAGATCAGTAGAT | CTAGCAAATTTCCTCTAC |
| JN543470_Bohadschia_sp: | 0.82 CTTATCTAGAAAAATTAGCCC | ATGCTGGGGATCAGTAGAC | CTTGCACATTTCCTTACA |
| JN543469_Bohadschia_sp: | 0.81 CTTATCTAGAAAAATTAGCCC | ATGCTGGGGATCAGTAGAC | CTTGCACATTTCCTTACA |
| NZEC0177508_Staurocucumis_liouvillei: | 0.81 TTATCAAGAAAATTGCTC | ATGCTGTTGGCTCTGTTGAT | CTGGCTATATTTCCTCTCA |
| HM196670_Staurocucumis_liouvillei: | 0.81 TTGTCAGAAAATTAGCCC | ATGCCGGAGGCTCTGTTGAC | CTAGCCATATTTCCTCTCA |
| JX683879_Bohadschia_koellikeri: | 0.81 CTTATCTAGAAAAATTAGCCC | ATGCTGGGGATCAGTAGAC | CTTGCACATTTCCTTACA |
| HM196671_Staurocucumis_liouvillei: | 0.80 TTATCAAGAAAATTGCGAC | ATGCTGGGGCTCCTGGTAT | CTAGCCATATTTCCTCTCA |
| NZEC017508_Araeosoma: | 0.80 ATTATCAAGAAAATTGCTC | ATGCTGGGGATCAGTAGAC | CTAGCCATCTTCATTTACA |
| AY262940_Echinometra_sp_C: | 0.80 ACTATCAAGAAAATTGCCC | ATGCTGGGGATCTGTTGAC | CTAGCAAATTTCCTCTCA |
| KF142181_Pentamera_calicera: | 0.80 TTATCAAGAAAATTGCCC | ATGCGAGGAGATCAGTTGAT | CTTGCCTATTTCCTCTACA |
| GU480560_Eucidaris_meturaria: | 0.80 TCATTCAGAAAATTGCTC | ACGCCGGAGGCTCTGTTGAT | CTAGCTATCTTCCTCTACA |
| JX544968_Phryella_fragilis: | 0.80 CTTATCAAGAAAATTAGCCC | ACCGAGGAGATCAGTTGAC | TTGGCCATATTTCCTCTCA |
| NZEC019508_Araeosoma: | 0.80 ATTATCAAGAAAATTGCTC | ATGCTGGGGATCAGTAGAC | CTAGCCATCTTCATTTACA |
| NZEC021508_Araeosoma: | 0.80 ATTATCAAGAAAATTGCTC | ATGCTGGGGATCAGTAGAC | CTAGCCATCTTCATTTACA |
| FJ971395_Holothuria_scabra: | 0.80 ACTATCAAGAAAATTGCCC | ATGCCGGAGGATCTGTTGAC | CTAGCCATTTCTCACTACA |
| GU480561_Echinometra_mathaei: | 0.80 ACTATCAAGAAAATTGCCC | ATGCTGGGGATCAGTAGAC | CTAGCAAATTTCCTCTCA |
| HM196724_Psolidium_gainsi: | 0.80 TCATTCAGAAAATTGCTC | ACCGAGGAGATCAGTTGAT | CTTGCCTATTTCCTCTACA |
| NZEC057508_Araeosoma: | 0.80 ATTATCAAGAAAATTGCTC | ATGCTGGGGATCAGTAGAC | CTAGCCATCTTCATTTACA |
| JX544962_Phryella_fragilis: | 0.80 ATTATCAAGAAAATTAGCCC | ACCGAGGAGATCAGTTGAC | TTGGCCATATTTCCTCTCA |
| GU480572_Holothuria_fuscopunctata: | 0.80 ACTATCAAGAAAATTGCCC | ATGCCGGAGGATCTGTTGAC | CTAGCCATTTCTCACTACA |

| | |
|--|---|
| +Echinodermata (phylum) 98% | - |
| +Echinidea (class) 8% | |
| +Echinida (order) 3% | |
| +Echinometridae (family) 3% | |
| +Echinometra (genus) 3% | |
| +Clypeasteroidea (order) 2% | |
| +Clypeasteroidea_incertaesedis (family) 2% | |
| +Clypeasteroidea_incertaesedis_incertaesedis (genus) 2% | |
| +Echinothurioida (order) 1% | |
| +Echinothuriidae (family) 1% | |
| +Araeosoma (genus) 1% | |
| +Cidaroida (order) 2% | |
| +Cidaroida_incertaesedis (family) ~0% | |
| +Cidaroida_incertaesedis_incertaesedis (genus) ~0% | |
| +Cidaridae (family) 2% | |
| +Eucidaris (genus) 2% | |
| +Holothuroidea (class) 62% | |
| +Aspidochirotiida (order) 45% | |
| +Holothuriidae (family) 45% | |
| +Holothuria (genus) ~0% | |
| +Actinopyga (genus) 2% | |
| +Bohadschia (genus) 38% | |
| +Elasipodida (order) 3% | |
| +Elpidiidae (family) 3% | |
| +Piniegona (genus) 3% | |
| +Dendrochirotida (order) 14% | |
| +Phyllophoridae (family) 6% | |
| +Pentameria (genus) 3% | |
| +Phyrella (genus) 1% | |
| +Dendrochirotida_incertaesedis (family) ~0% | |
| +Dendrochirotida_incertaesedis_incertaesedis (genus) ~0% | |
| +Cucumariidae (family) 4% | |
| +Staurocucumis (genus) 4% | |
| +Psolididae (family) 3% | |
| +Psolus (genus) 1% | |
| +Psolidium (genus) 2% | |
| +Ophiuroidea (class) 18% | |
| +Ophiurida (order) 12% | |
| +Ophiacanthidae (family) 10% | |
| +Ophiacantha (genus) 10% | |
| +Ophiotrichidae (family) 1% | |
| +Ophiomaza (genus) 1% | |
| +Euryalida (order) 5% | |
| +Euryalida_incertaesedis (family) 5% | |
| +Euryalida_incertaesedis_incertaesedis (genus) 5% | |
| +Asteroidea (class) ~0% | |

REMOTE vs LOCAL DATABASES

BLAST against Genbank

- Slow!
- Bad taxonomy!
- Misidentified entries!



BOLD identification

- Good, curated, taxonomy!
- Only relevant sequences
- Slow for too many seqs
- Only for conventional barcodes!



ECOTAG using a local database

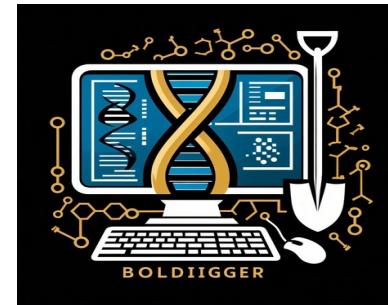
- Usually fast (although the speed depends on the reference database size!)
- Errors and non-relevant sequences may be removed from the database
- Custom sequences may be added immediately
- Can be used for non-conventional markers and novel primer sets



OTHER TAXONOMIC ASSIGNMENT ALGORITHMS

BOLDigger / Boldigger II (Buchner & Leese 2020)

<https://github.com/DominikBuchner/BOLDigger2>



The BOLD identification tool allows only batches of 100 sequences that can be identified in one run.

BOLD's API does not grant access to private and early release data.

BOLDigger aims to solve this problem. It gives automated access to the identification engine and can also be used to download additional metadata for each sequence as well as helping to choose the top hit from the returned results.

mkLTG (Meglécz 2024)

Faster than ecotag, also based in the LCA approach to phylogenetic assignment, however with arbitrary thresholds for taxonomic ranks

It comes with an automatically-generated database: COInr

> Biol Futur. 2023 Dec;74(4):369-375. doi: 10.1007/s42977-024-00201-x. Epub 2024 Feb 1.

mkLTG: a command-line tool for taxonomic assignment of metabarcoding sequences using variable identity thresholds

Emese Meglécz ¹

<https://github.com/meglecz/mkLTG>

<https://github.com/meglecz/mkCOInr>

OTHER TAXONOMIC ASSIGNMENT ALGORITHMS

MOLECULAR ECOLOGY
RESOURCES

WILEY

Assessment of current taxonomic assignment strategies for metabarcoding eukaryotes

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

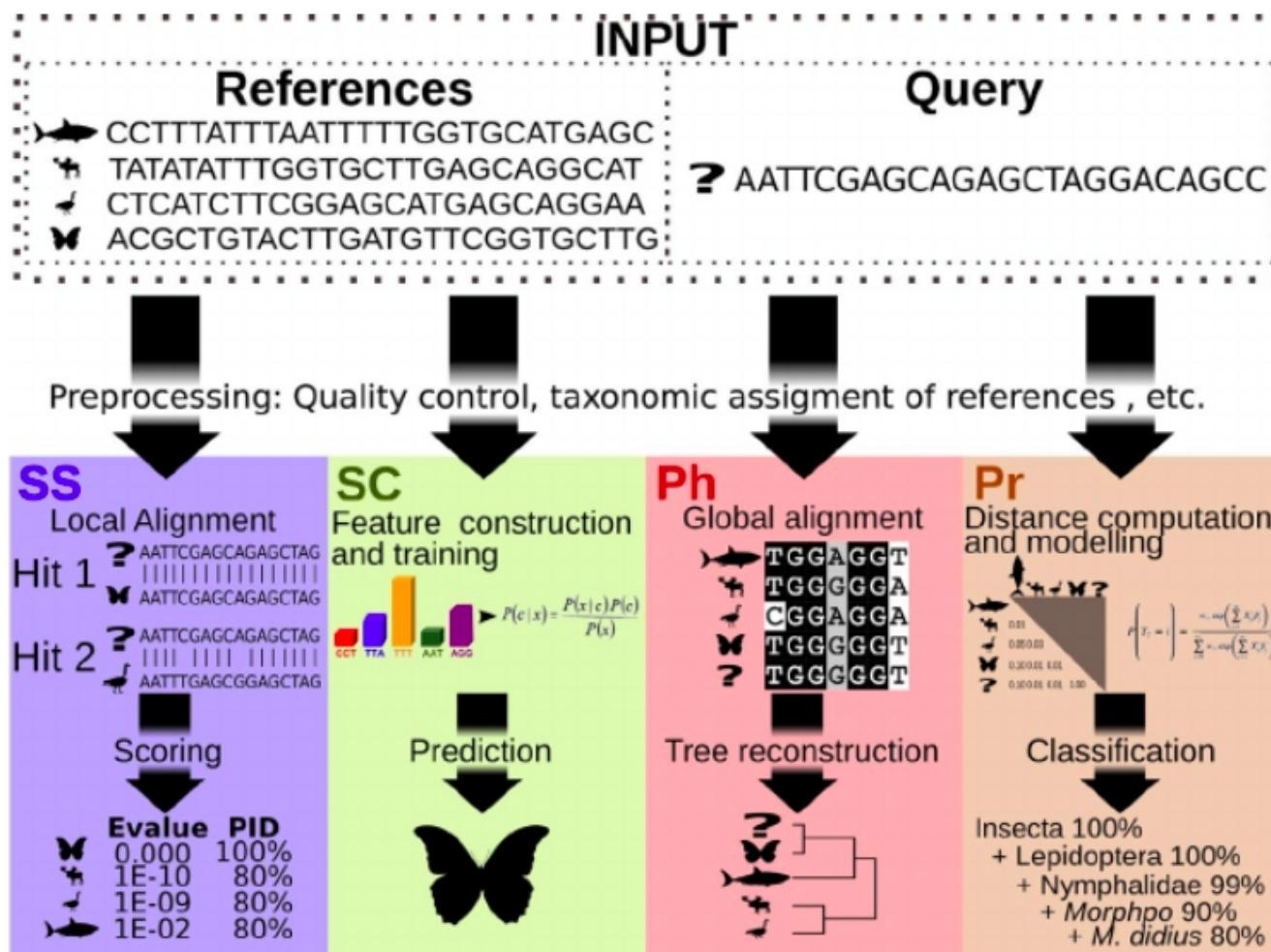


FIGURE 1 Overview of methodological approaches. Sequence similarity (SS) methods use local alignments to search for similarity between each query and the reference sequences. Sequence composition (SC) methods are trained by computing a k-mer frequency profile for each reference sequence, and then matching each query to this profile. Phylogenetic (Ph) methods use global alignments including (or placing) the query in a phylogenetic tree. Probabilistic (Pr) methods use a distance metric and then perform a hierarchical multinomial regression to estimate the certainty in the classification of each query at each taxonomic rank

BIOINFORMATICS PIPELINE: REFINING THE FINAL DATASET

- Removal of false positives from tag-switching
- Taxonomic clustering (cluster MOTUs assigned to same taxon)
- Removal of sequences derived from pseudogenes: LULU
- Blank corrections
- Removal of sequences derived from contaminations
- Rarefaction and re-sampling techniques for controlling variability in total reads

COMMONLY USED METABARCODING PIPELINES

- QIIME / QIIME2: <https://qiime2.org/>
- Mothur: <https://www.mothur.org/>
- OBITools: <https://metabarcoding.org/obitools>
- JAMP: <https://github.com/VascoElbrecht/JAMP>
- USEARCH: <https://www.drive5.com/usearch/>
- VSEARCH: <https://github.com/torognes/vsearch>
- DADA2: <https://benjneb.github.io/dada2/>
- APSCALE: <https://github.com/DominikBuchner/apscale>
- MJOLNIR: <https://github.com/uit-metabarcoding/MJOLNIR>

SOME CURRENT APPLICATIONS

STUDIES ON INDIVIDUALS

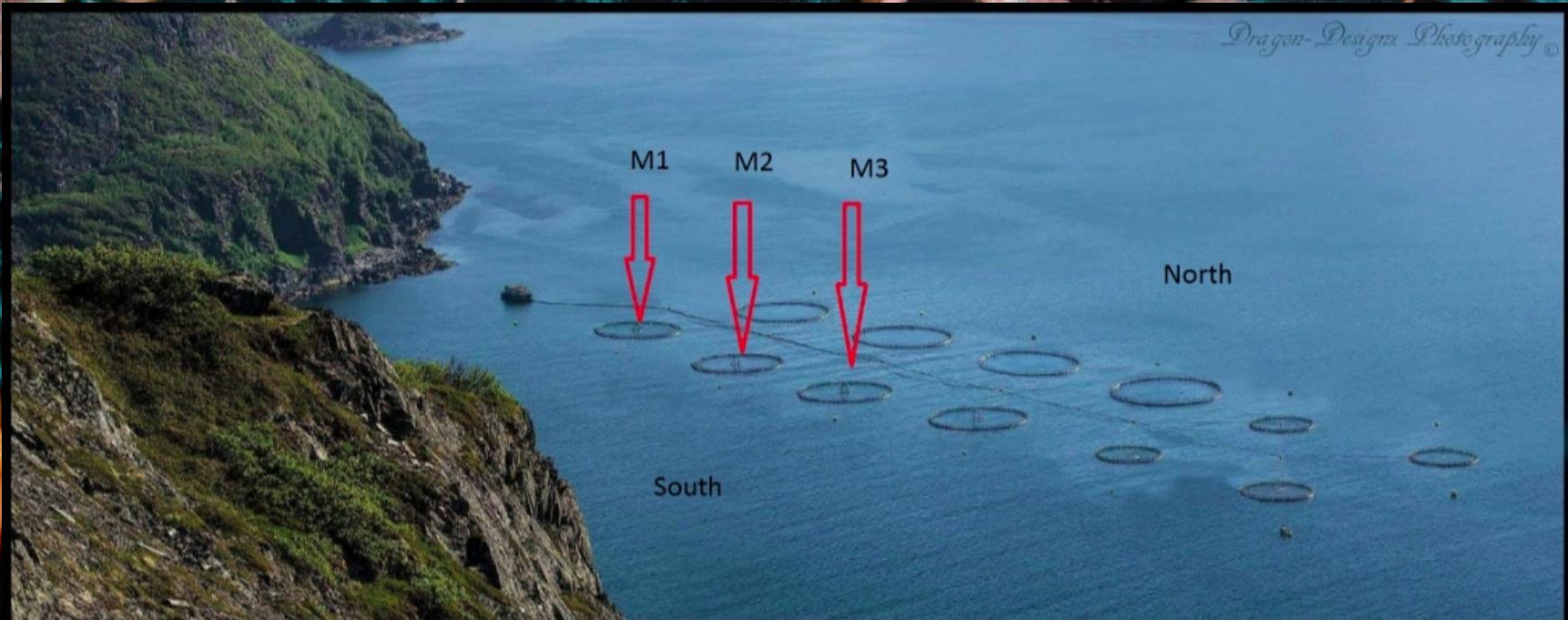
- Trophic relationships: gut contents / analyses of faeces
- Symbiotics: host/pathogen/symbiont interactions

STUDIES ON COMMUNITIES: SPATIAL PATTERNS

- Molecular characterization of biomes
- Population metagenomics / Metaphylogeography

STUDIES ON COMMUNITIES: TEMPORAL PATTERNS

- Community dynamics / annual cycles
- Settlement events detection / reproductive cycles
- Long-term shifts in community structure (including genetic structure)
- Distribution shifts / assessment of effects of Climate Change
- Detection / monitoring of introduced species
- Studies on palaeoecological records: palaeometabarcoding



**AND EXAMPLE USING THE
MJOLNIR PIPELINE**