



Metabarcoding : the main steps - Part 1

Clarisse Lemonnier

The INRAE logo is positioned at the bottom left of the slide. It consists of the letters "INRAE" in a bold, teal, sans-serif font. The letter "E" is stylized with a circular element at its base. The logo is partially overlaid by a large, abstract graphic of overlapping hexagons in various shades of green and teal that occupies the left side of the slide.

INRAE



Summary

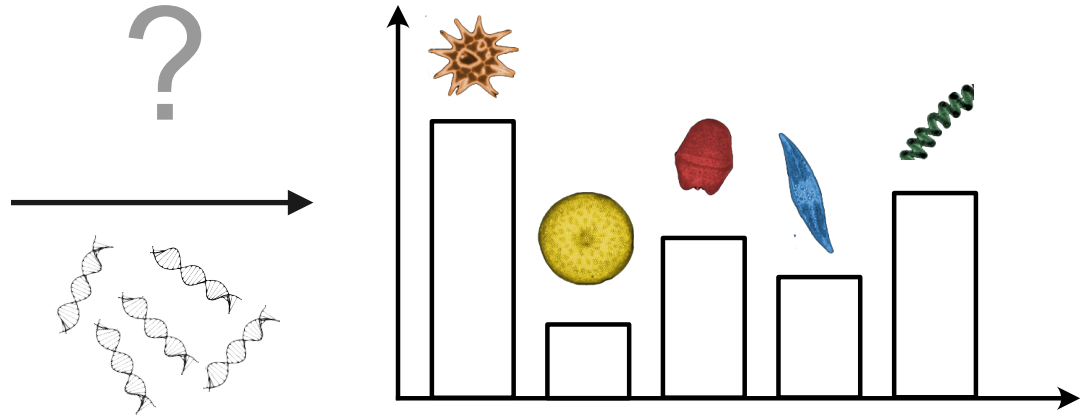
The main steps of metabarcoding

Selection of the barcode

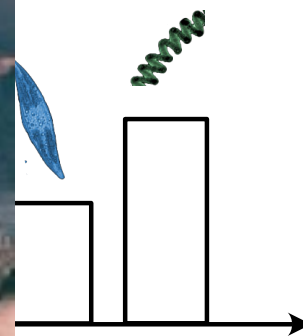
Reference databases



Metabarcoding steps



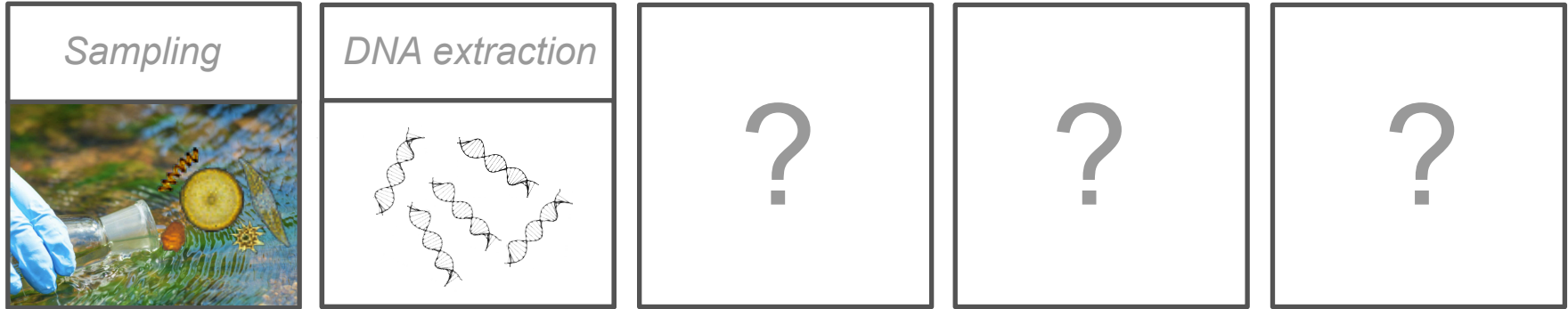
Metabarcoding steps



Metabarcoding steps



Metabarcoding steps



Metabarcoding steps

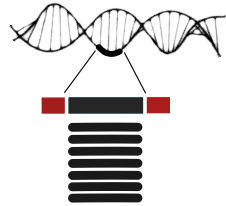
Sampling



DNA extraction



*Barcode
amplification*



?

?

Metabarcoding steps

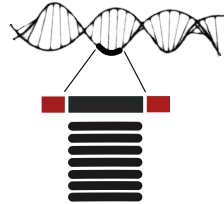
Sampling



DNA extraction



*Barcode
amplification*



Sequencing

```
ATCGCTTTGGACCT  
ATCGAATTGGAACA  
ATCGCTTTGGACCT  
ATCGAATTGGAACA  
ATCGCTTTGGACCT  
ATCGAATTGGAACA  
ATCGCTTTGGACCT  
ATCGAATTGGAACA
```



?

Metabarcoding steps

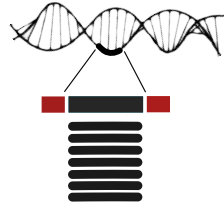
Sampling



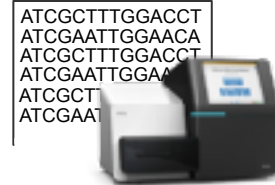
DNA extraction



Barcode amplification

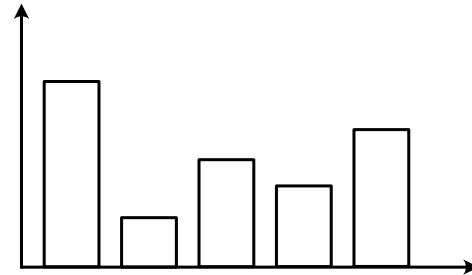


Sequencing

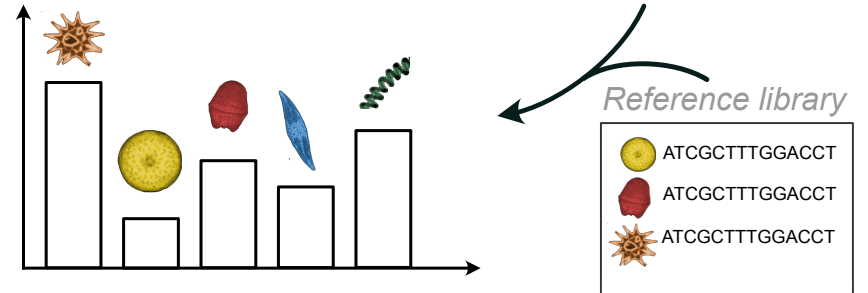
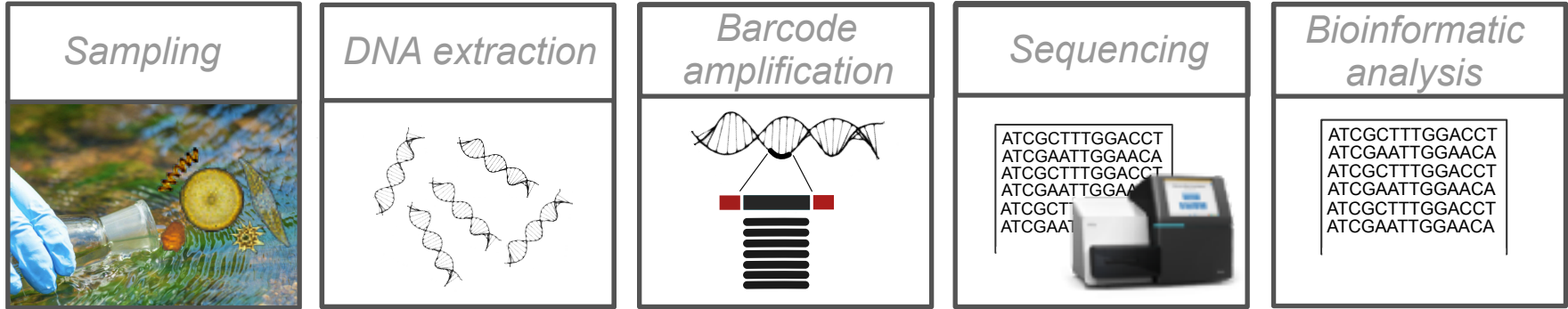


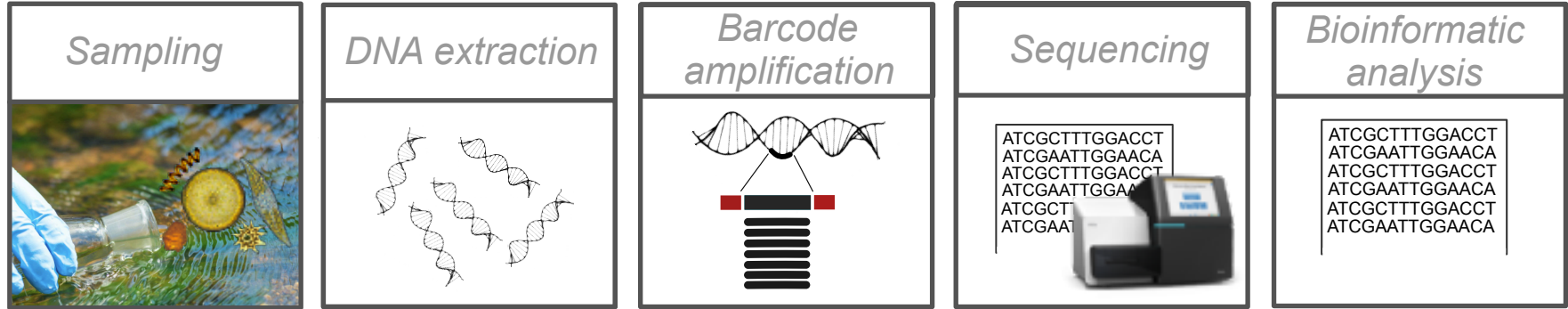
Bioinformatic analysis

```
ATCGCTTTGGACCT
ATCGAATTGGAACA
ATCGCTTTGGACCT
ATCGAATTGGAACA
ATCGCTTTGGACCT
ATCGAATTGGAACA
ATCGAATTGGAACA
```

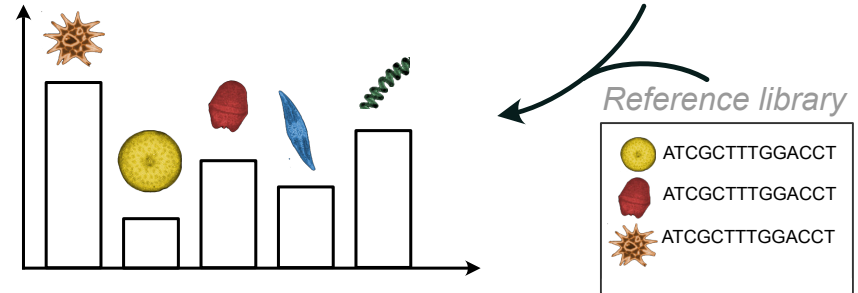


Metabarcoding steps

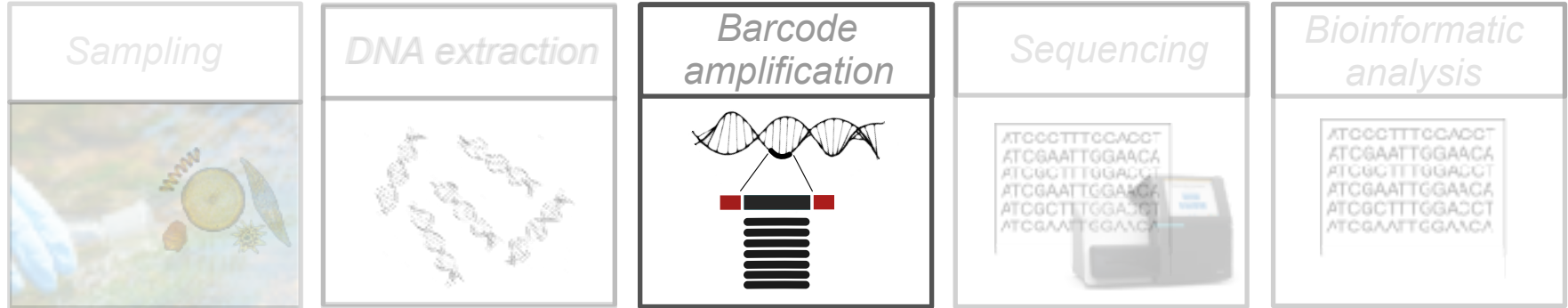




An informed choice at each step is required to get a reliable result at the end

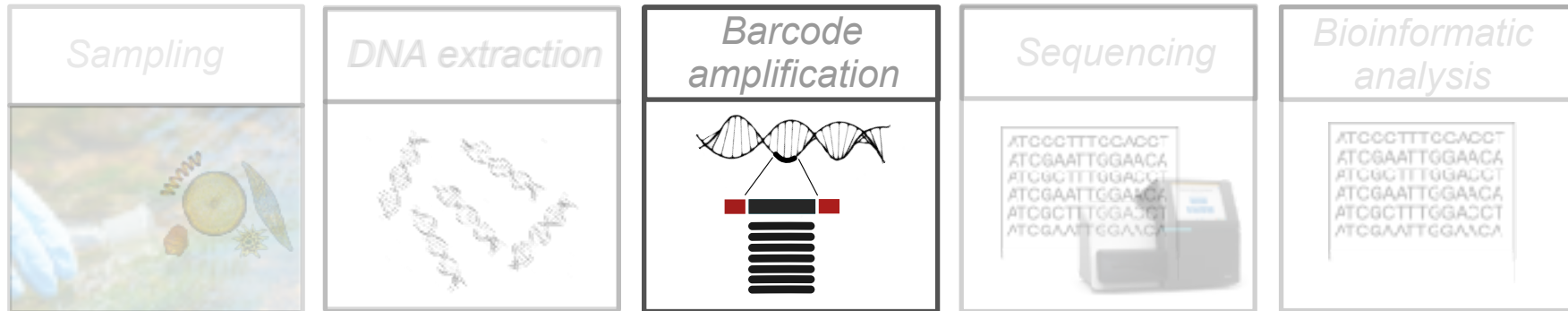


Metabarcoding steps



Barcode selection

Metabarcoding steps

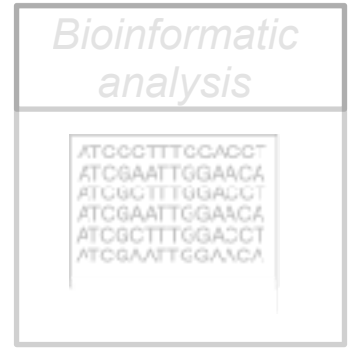
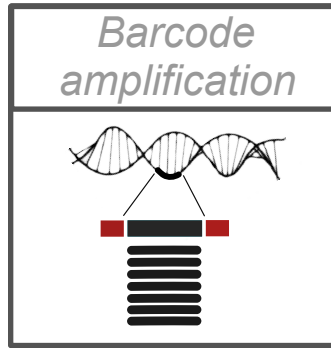


Recovery of all species



Barcode selection

Metabarcoding steps

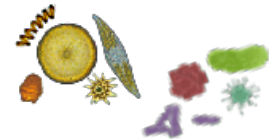


Recovery of all
species

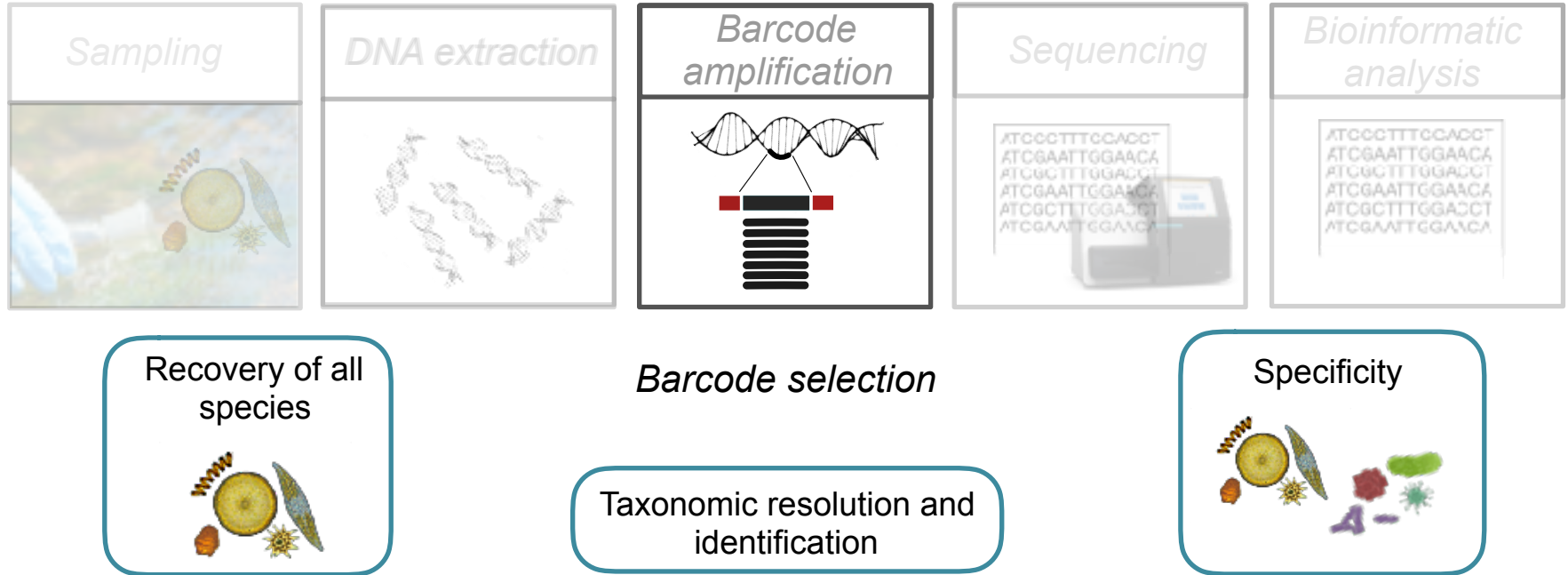


Barcode selection

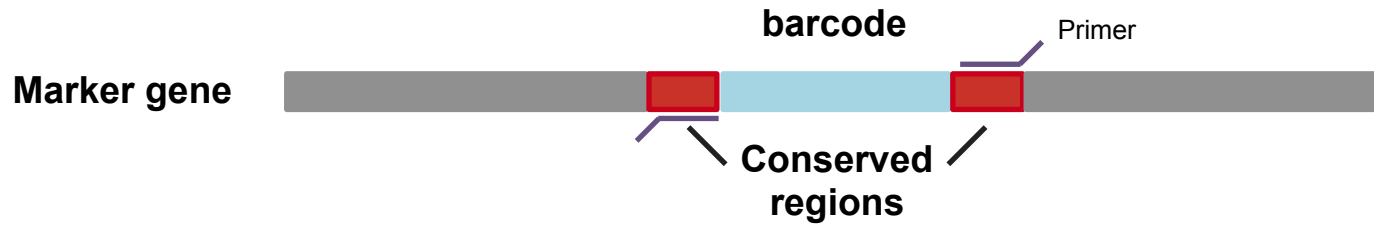
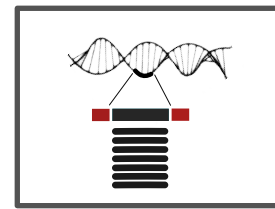
Specificity

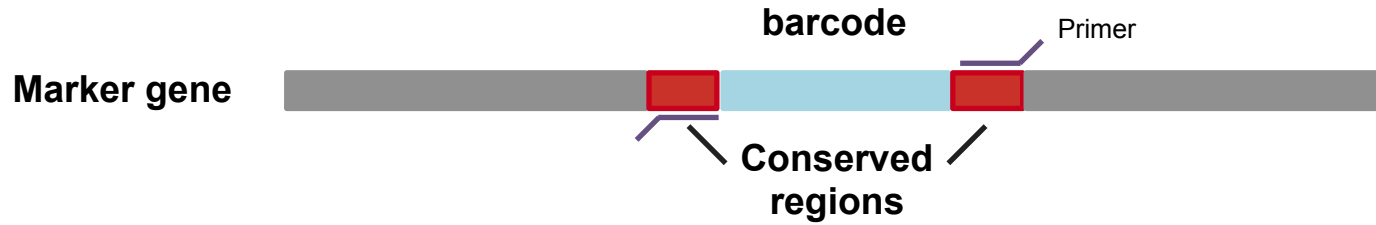
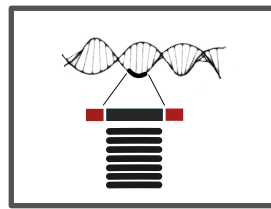


Metabarcoding steps



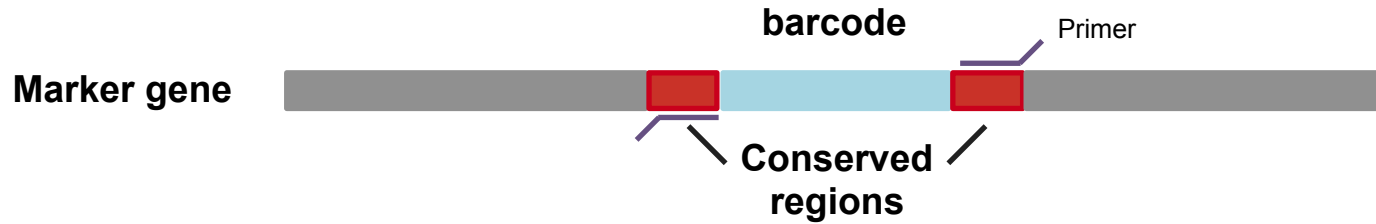
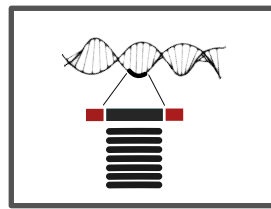
Barcode selection - a fundamental criteria





1 Universal

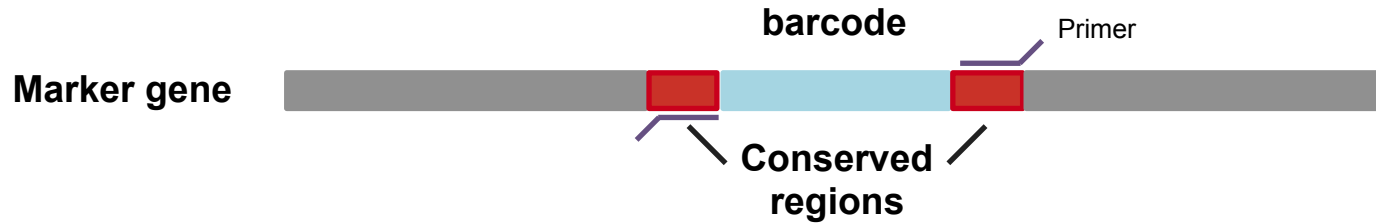
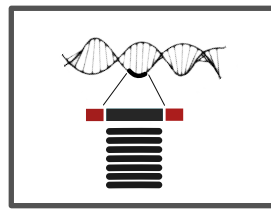




1 Universal

2 Conserved regions

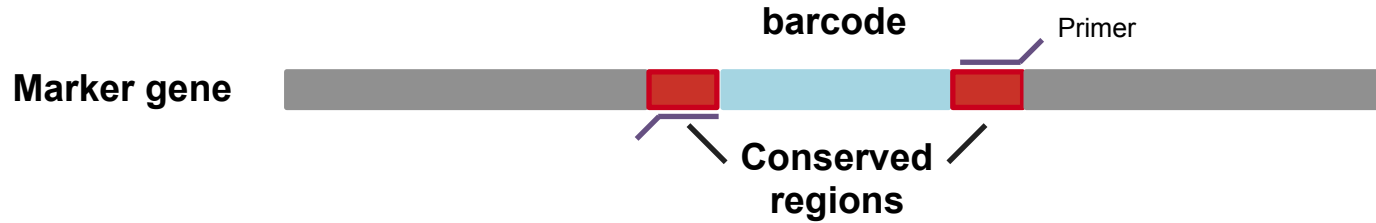
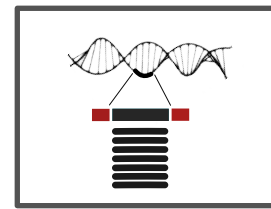




- 1 Universal
- 2 Conserved regions
- 3 Variable enough



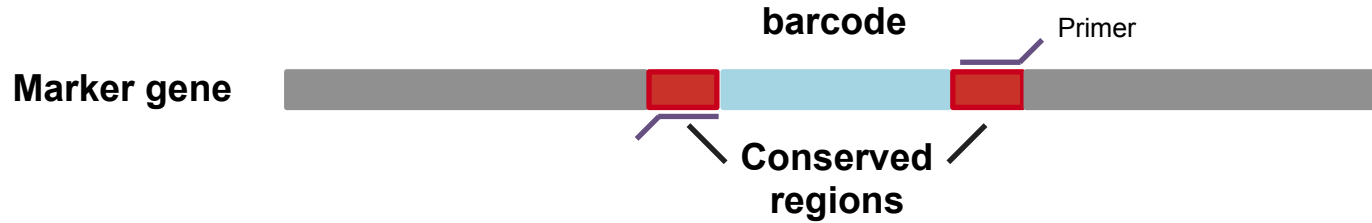
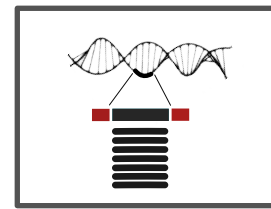
Barcode selection - a fundamental criteria



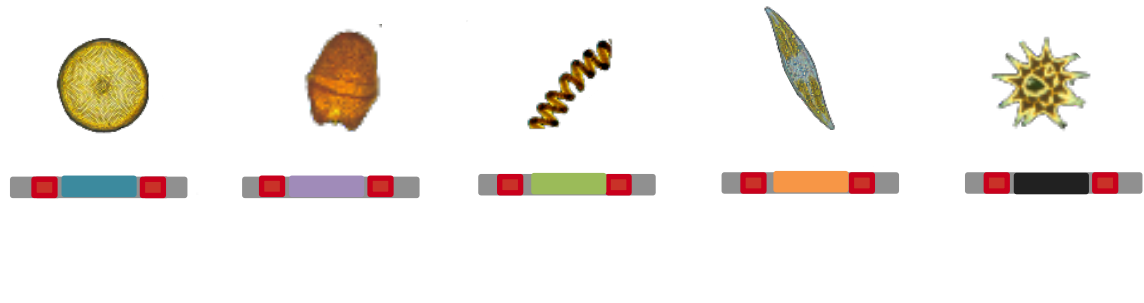
- 1 Universal
- 2 Conserved regions
- 3 Variable enough
- 4 Match sequencing technology size



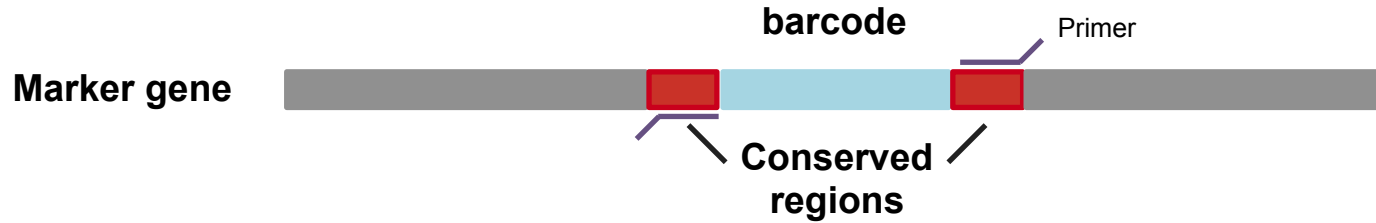
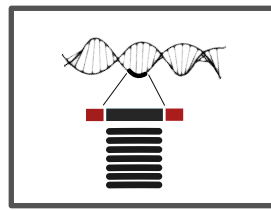
Barcode selection - a fundamental criteria



- 1 Universal
- 2 Conserved regions
- 3 Variable enough
- 4 Match sequencing technology size
- 5 Represented in reference libraries



Barcode selection - a fundamental criteria



1 Universal

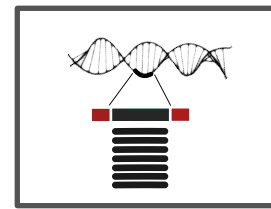
2 Conserved regions

3 Variable enough

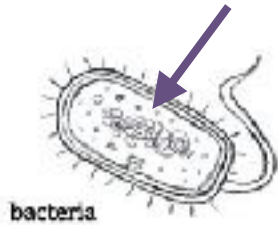
4 Match sequencing technology size

5 Represented in reference libraries

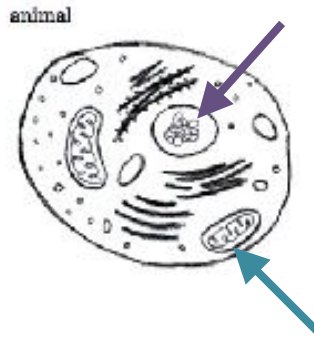
What are the classical marker genes used in metabarcoding?



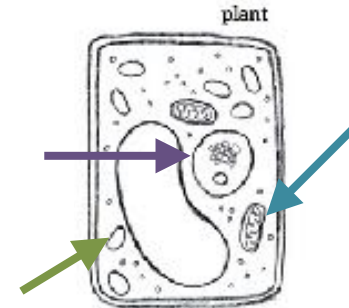
- A marker gene can be present in all organisms genomes



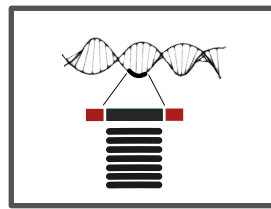
Genomic DNA



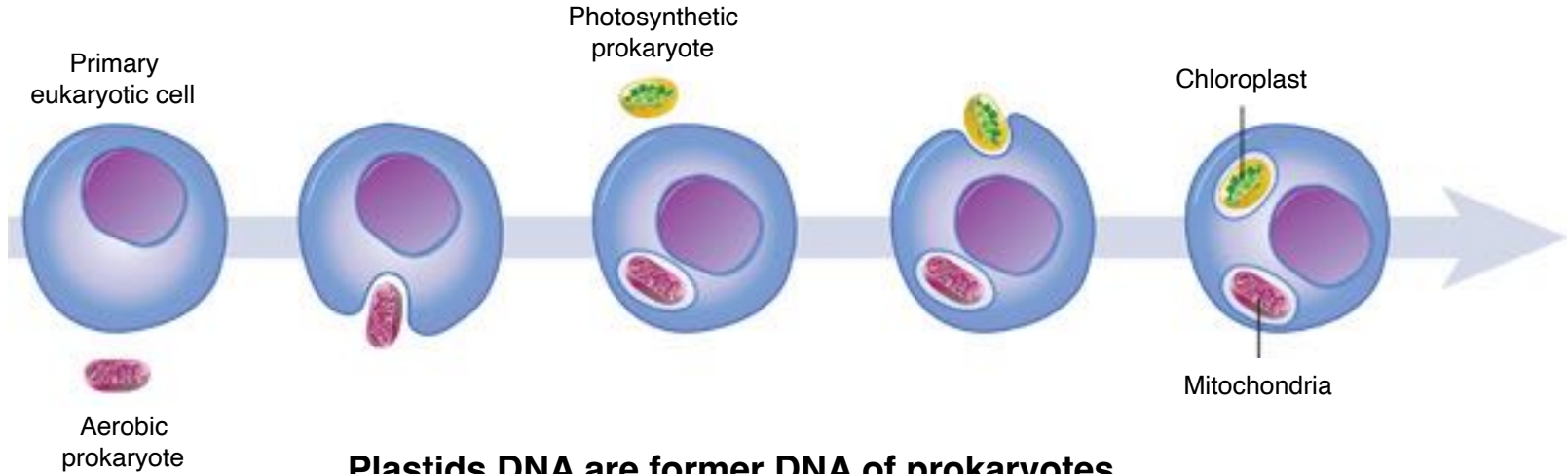
Mitochondrial DNA



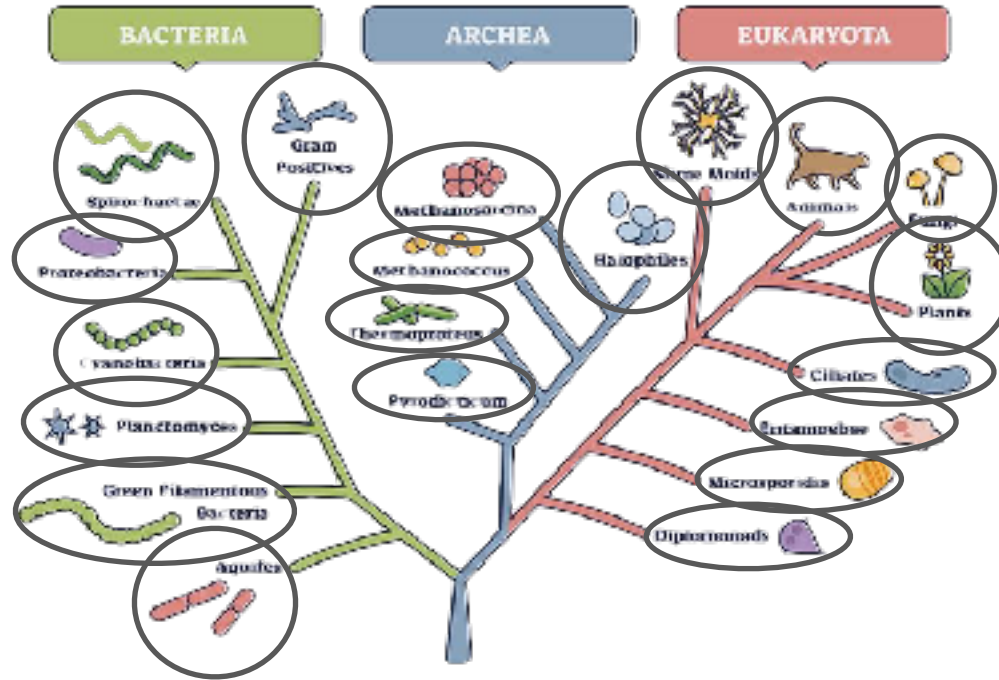
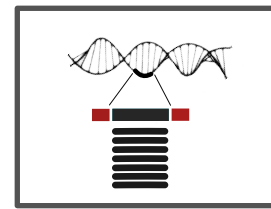
Chloroplast DNA



- A marker gene can be present in all organisms genomes



Classical marker genes - an overview



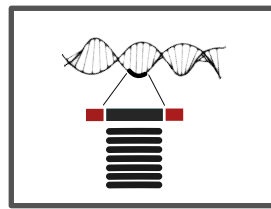
Ribosomal proteins

Nuclear

Mitochondrial

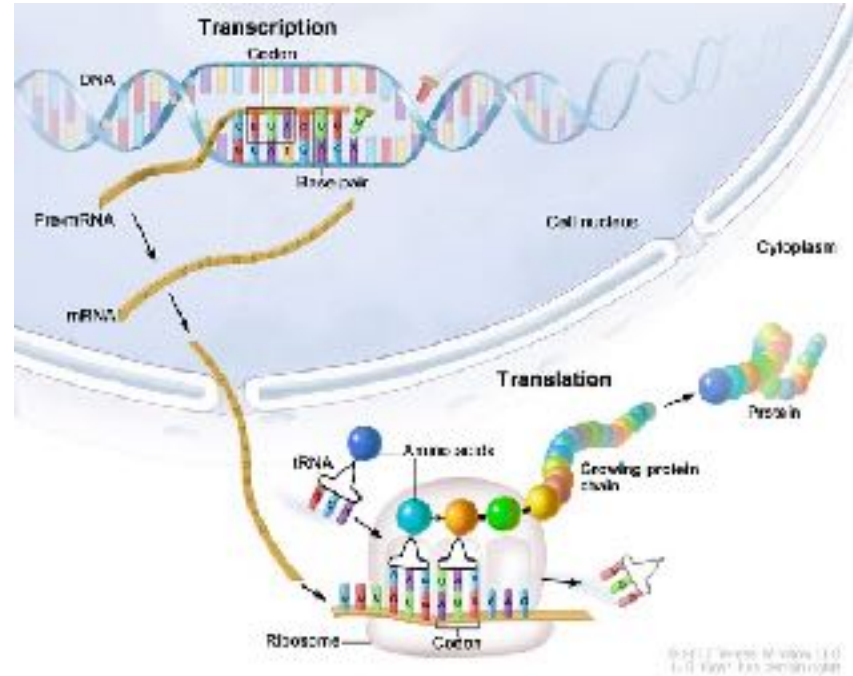
Chloroplast

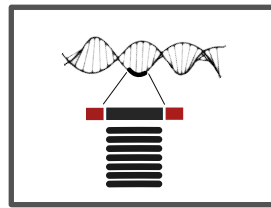




Ribosomal genes

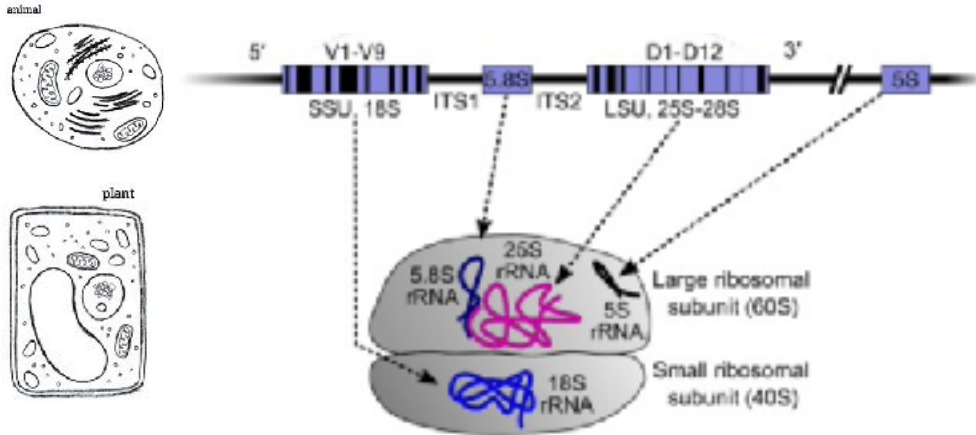
- Code for the large and small units of ribosomes which are responsible of the translation of mRNA into proteins
- Present in the nuclear DNA of **all living organisms**
- Present in the plastid DNA (mitochondria, chloroplasts)

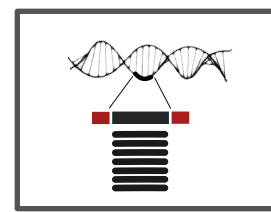




Ribosomal genes

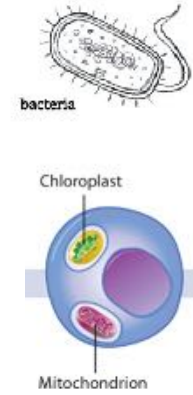
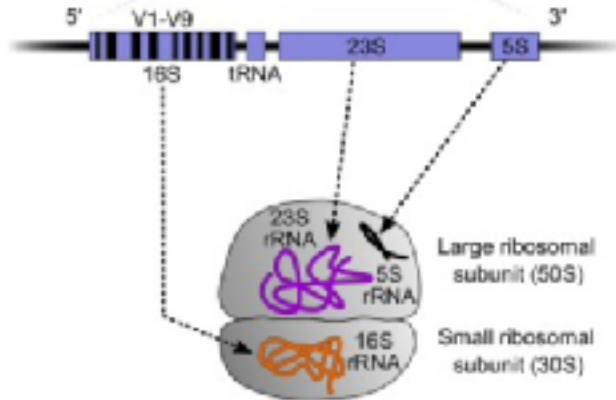
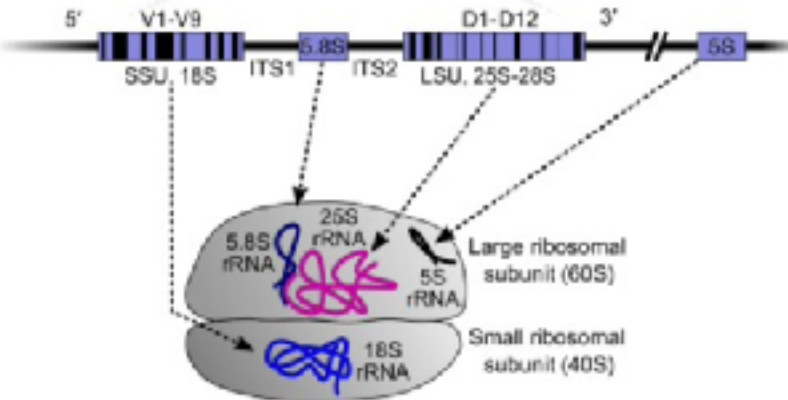
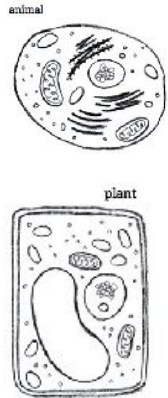
- Genes organized in an operon (*i.e.* clustered in the genome)



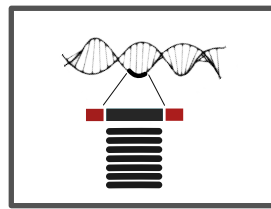


Ribosomal genes

● Genes organized in an operon (*i.e.* clustered in the genome)

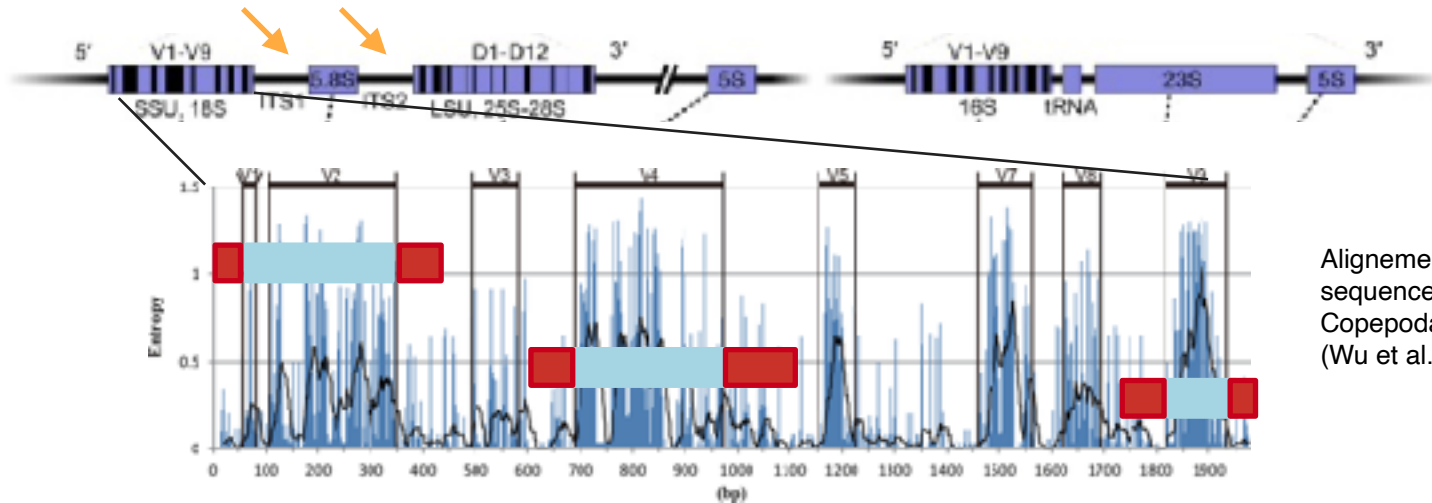


Trends in Microbiology

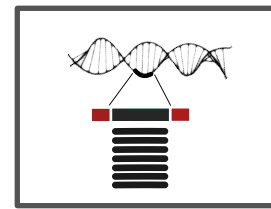


Ribosomal genes

- Alternance of hyper variable and conserved regions + Internal Transcribe spacer (ITS)

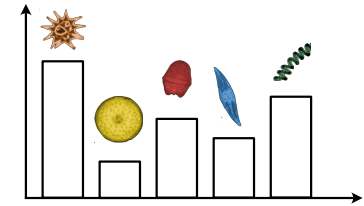
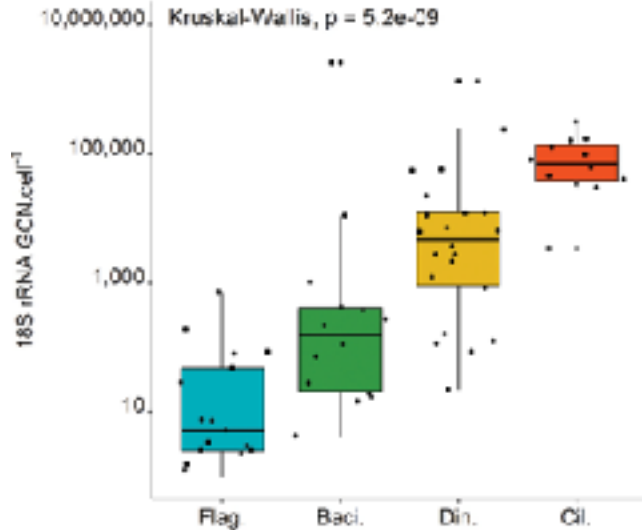


Alignment of 18S sequences from 192 Copepoda species (Wu et al., 2015)



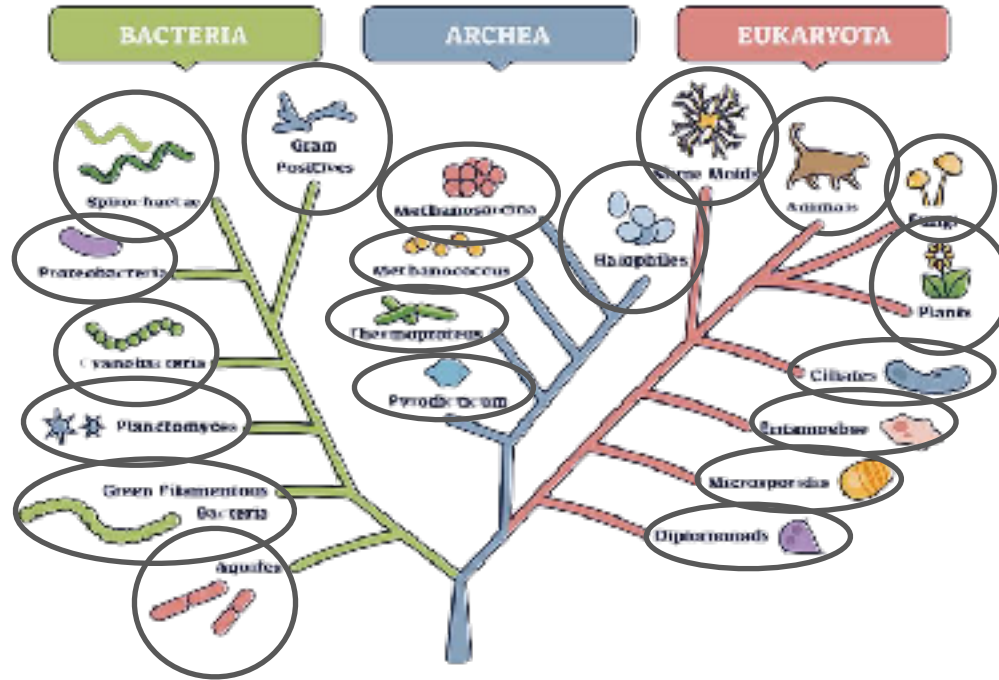
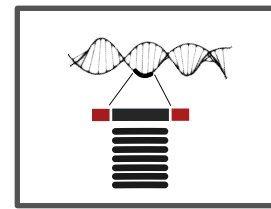
Various copy numbers between taxa

To meet the huge biosynthetic demand, eukaryotic cells contains hundreds to thousands of copies of ribosomal genes. (Kobayashi et al., 2011)



Influence on abundance

Classical marker genes - an overview



Ribosomal proteins

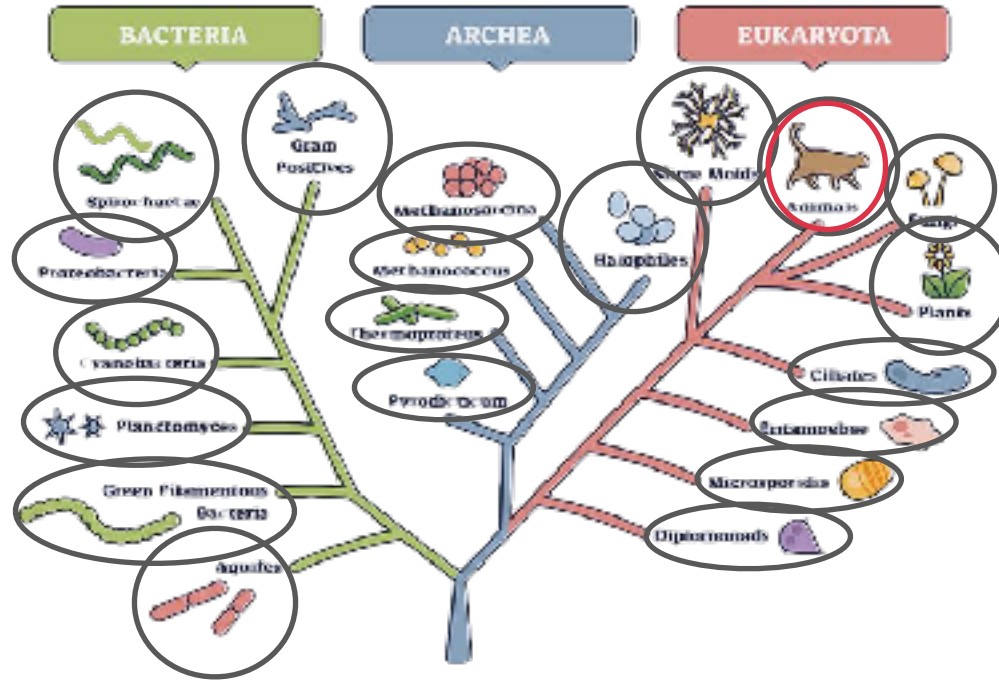
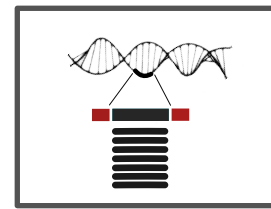
Nuclear

Mitochondrial

Chloroplast



Classical marker genes - an overview



Ribosomal proteins

Nuclear

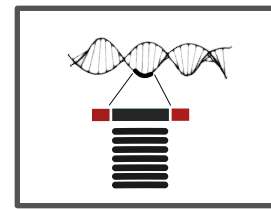
Mitochondrial

Chloroplast

COI

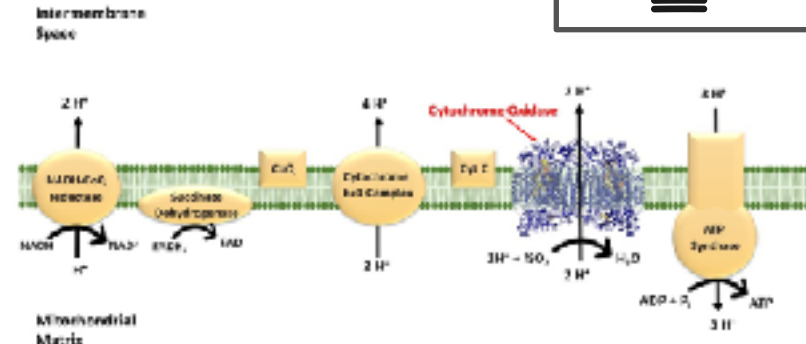
Mitochondrial

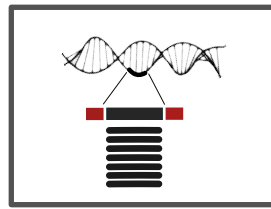




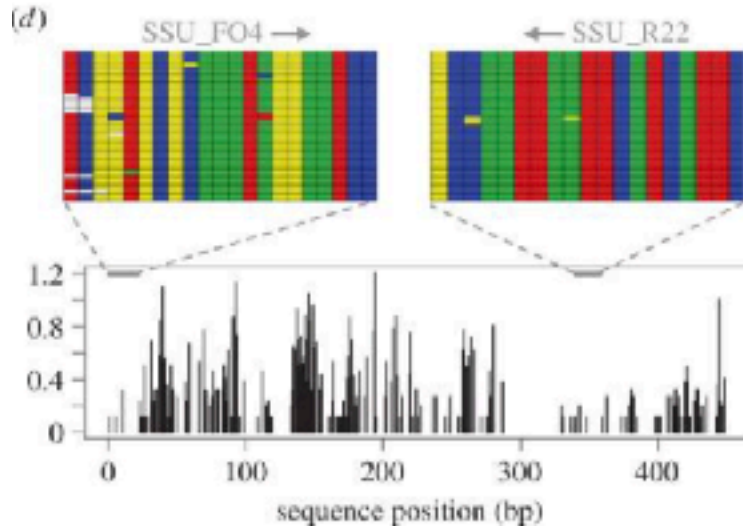
COI

- COI codes for the Cytochrome Oxidase subunit I. It is a key mitochondrial enzyme for respiration.
- Historically it is used for barcoding of animals. A lot of primers has been designed for various animals group.
- COI often provides a better taxonomic resolution than nuclear rRNA (18S).
- But contrary to ribosomal genes, it is complex to define universal primers.

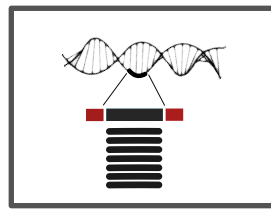




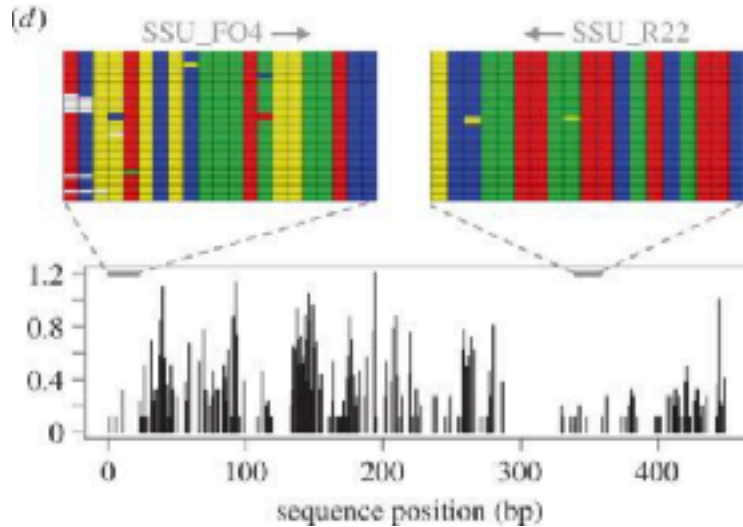
Nuclear 18S



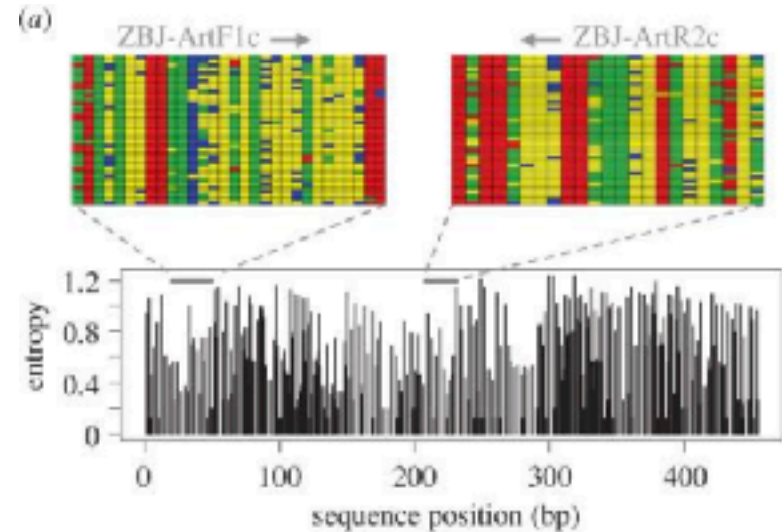
From more than 40 species of insects among 25 different orders (Deagle et al. 2014)



Nuclear 18S

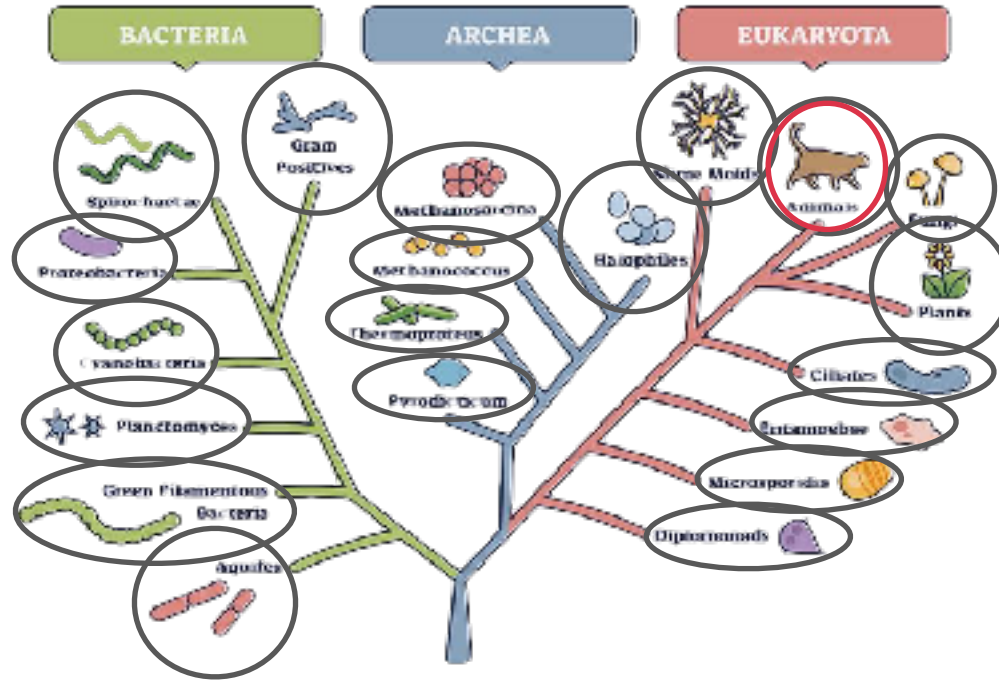
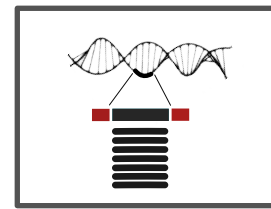


COI



From more than 40 species of insects among 25 different orders (Deagle et al. 2014)

Classical marker genes - an overview



Ribosomal proteins

Nuclear

Mitochondrial

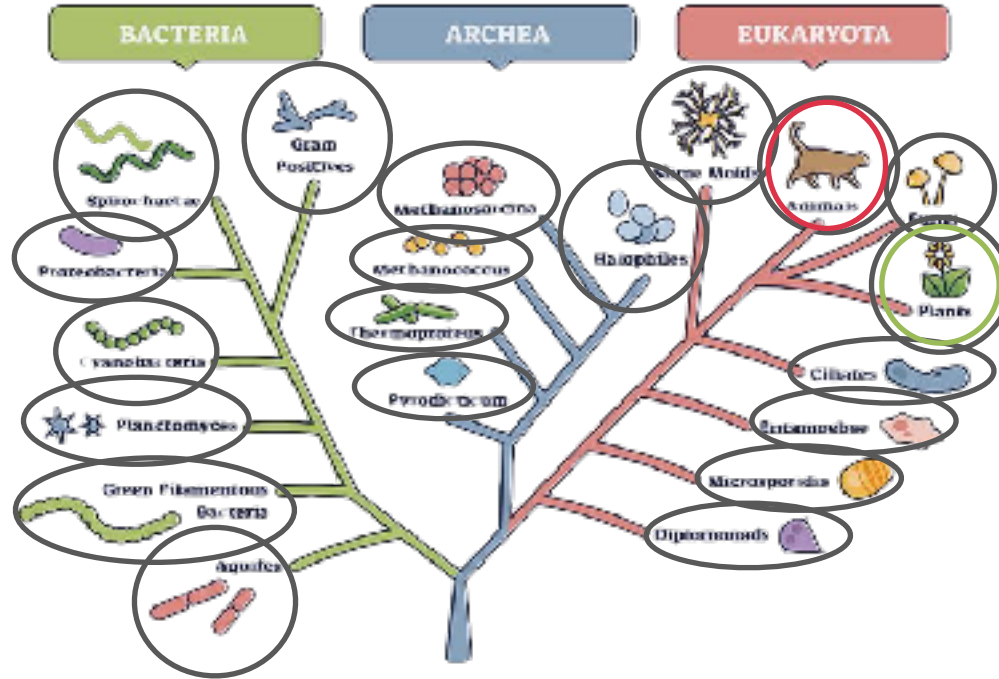
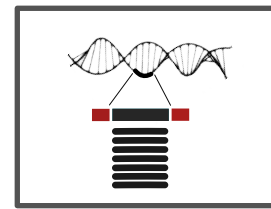
Chloroplast

COI

Mitochondrial



Classical marker genes - an overview



Ribosomal proteins

Nuclear

Mitochondrial

Chloroplast

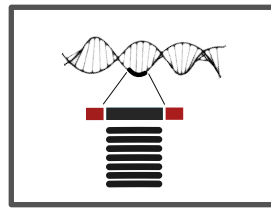
COI

Mitochondrial

RbcL/matK

Chloroplast





Rbcl/MatK - chloroplastic genes to target photosynthetic organisms

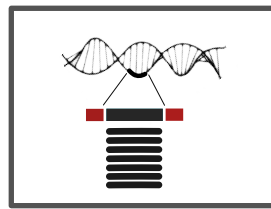
- For plant, low substitution rates of mitochondrial DNA - COI is not a good target
- Genes present in chloroplasts - target specifically photosynthetic organisms

MatK : plastid genes that codes for an intron maturase

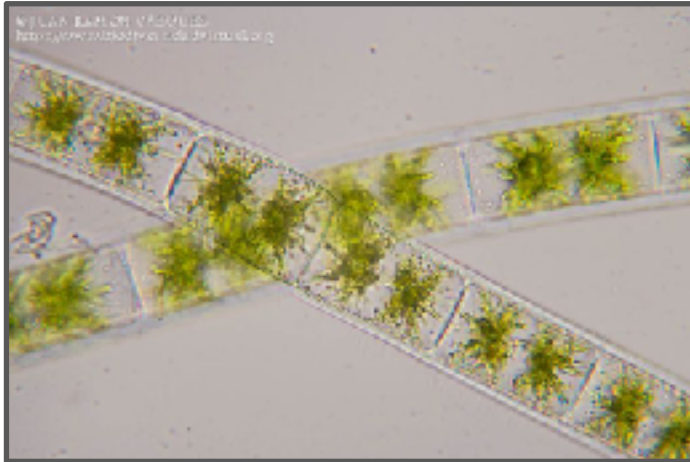
Rbcl : codes for the large subunit of RuBisCo, a key enzyme of photosynthesis.

Difficulties to get target all the phytoplanktonic diversity

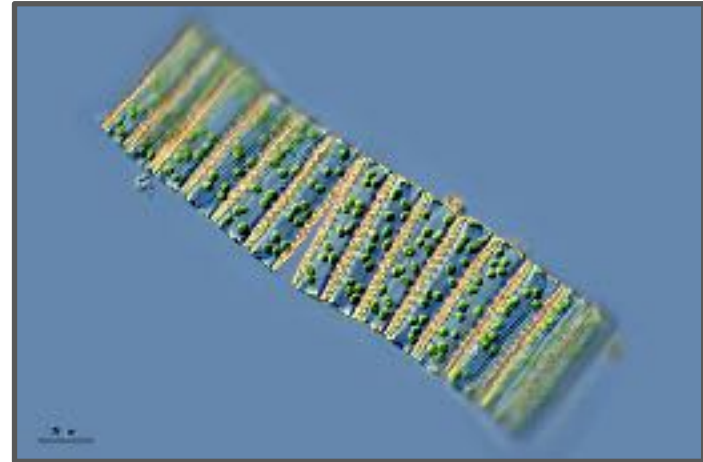
Ribosomal genes (16S)



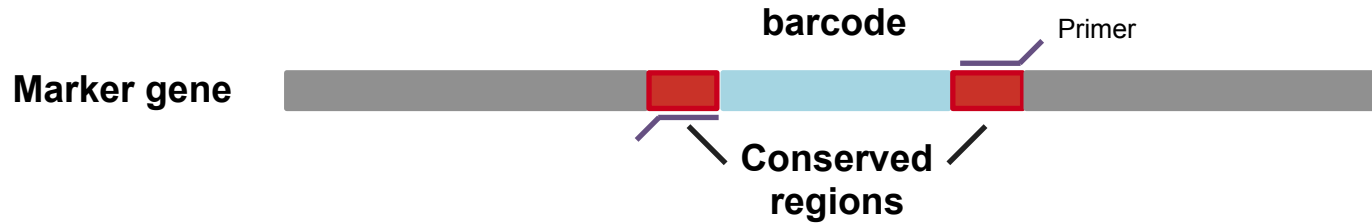
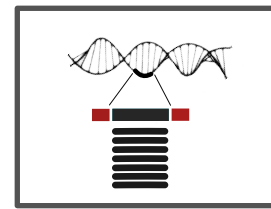
A marker gene in chloroplasts and mitochondria can be present in variable numbers in the cell



Zygnema sp. 1-2 chloroplasts/cell



Diatoma sp. > 10 chloroplasts/cell



1 Universal

2 Conserved regions

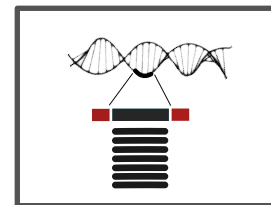
3 Variable enough

4 Match sequencing technology size

5 Represented in reference libraries

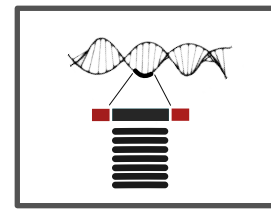
How primers are designed to hybridize only in the DNA of all targeted organisms?

Specificity



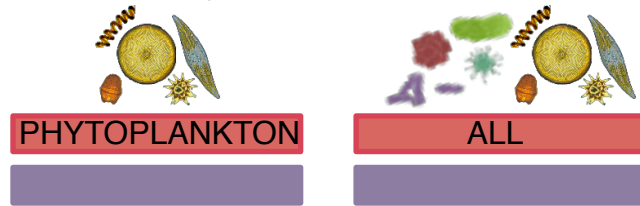
Conserved regions : specificity

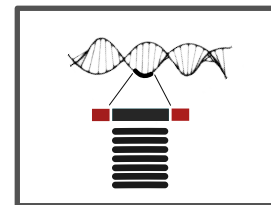




Conserved regions : specificity

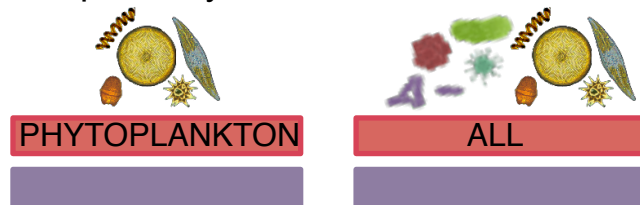
Primer specificity :



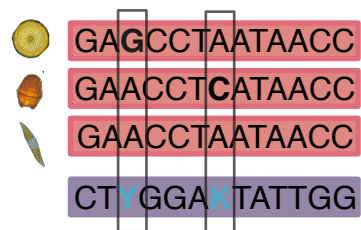


Conserved regions : specificity

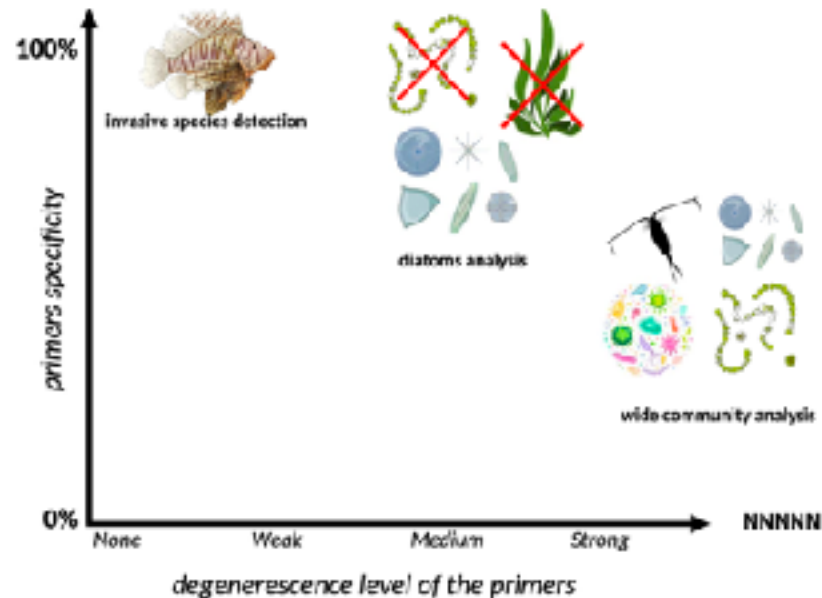
Primer specificity :



Primer degenerescence :

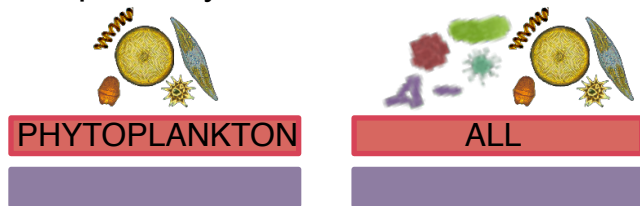


Code	Description
M	AC
R	AG
W	AT
S	CG
Y	CT
K	GT
V	ACG
H	ACT
D	AGT
B	CGT
N	ACGT

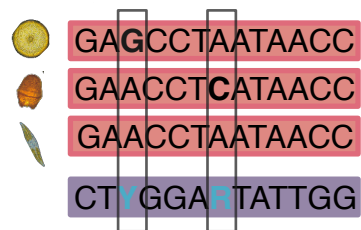


Conserved regions : specificity

Primer specificity :

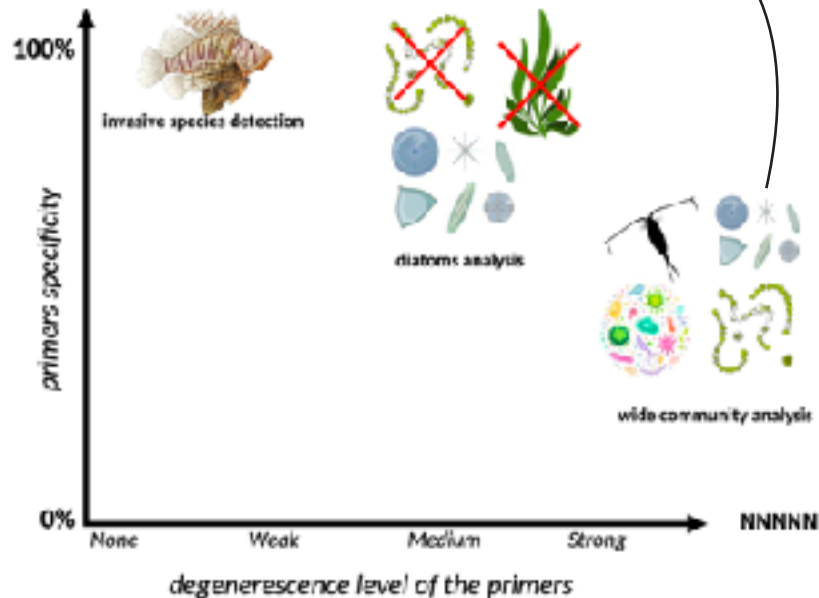
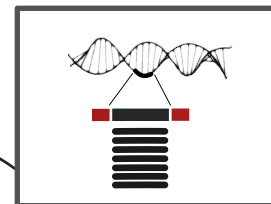


Primer degenerescence :

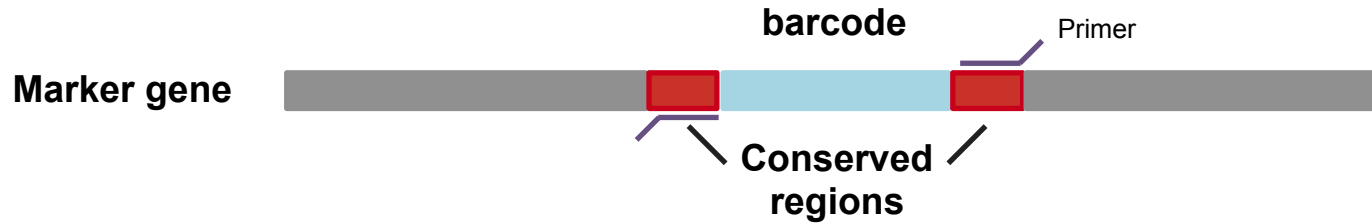
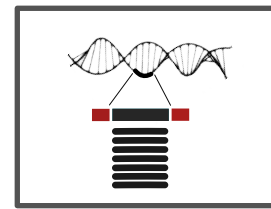


Code	Description
M	AC
R	AG
W	AT
S	CG
Y	CT
K	GT
V	ACG
H	ACT
D	AGT
B	CGT
N	ACGT

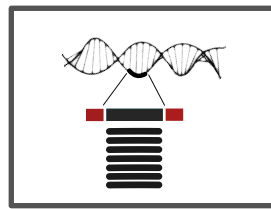
515F GTGYTCAGCMGCCGCGGTAA
926R CCGYCAATTYMTTTRAGTTT



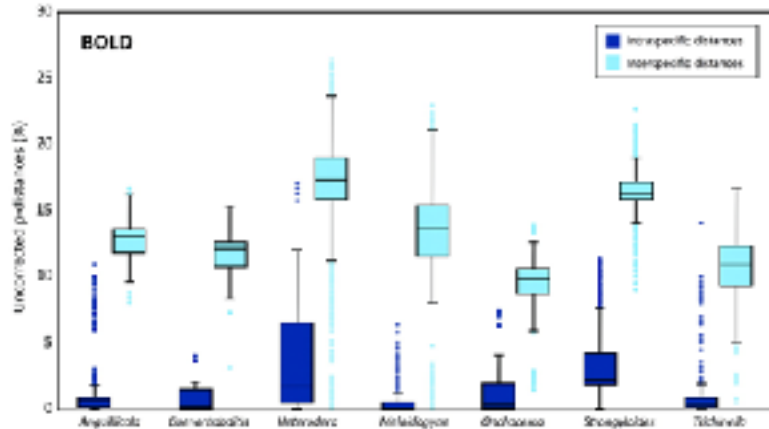
Barcode selection - a fundamental criteria



- 1 Universal
- 2 Conserved regions
- 3 Variable enough
Barcode-gap efficiency for a great diversity of organisms in metabarcoding
- 4 Match sequencing technology size
- 5 Represented in reference libraries

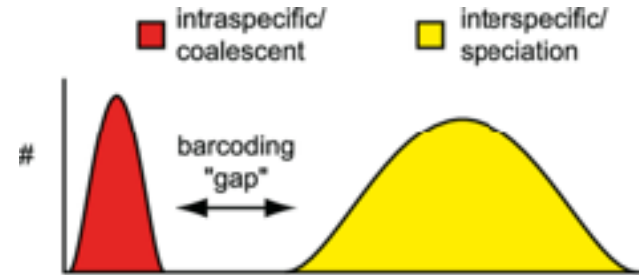


Depends on the taxonomic group



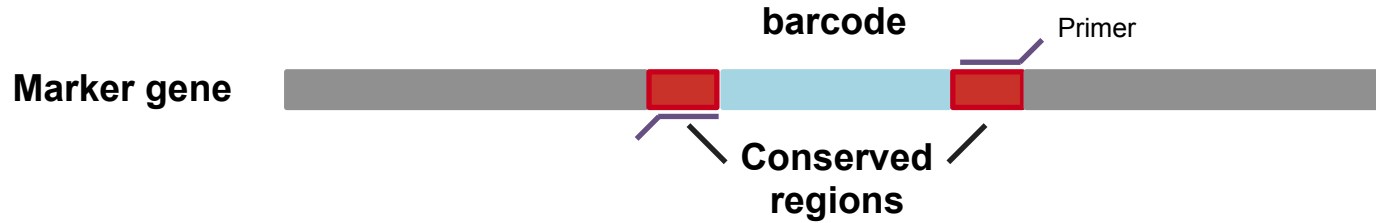
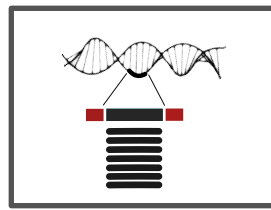
Tresoldi-Gonsaves et al., 2021

COI of various genera of worms (Nematoda)



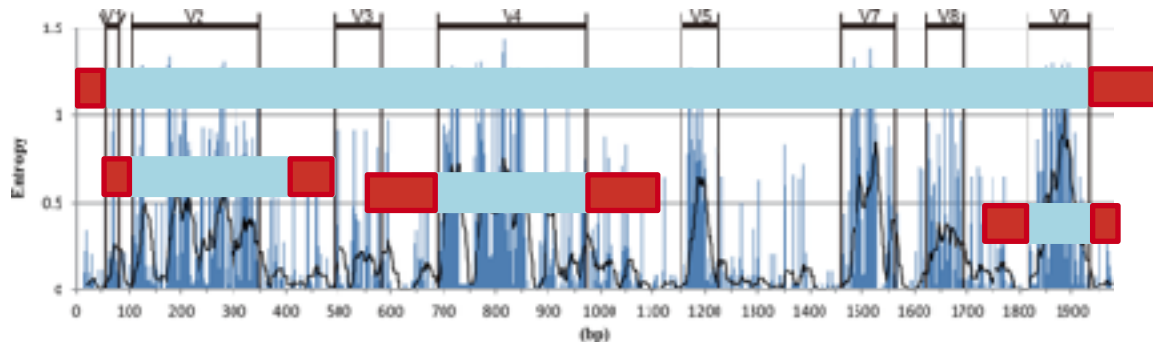
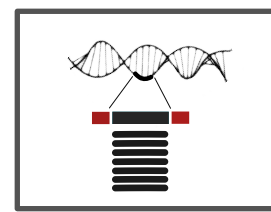
There is often a trade-off between capacity to target a large diversity and the barcode variability

Barcode selection - a fundamental criteria



- 1 Universal
- 2 Conserved regions
- 3 Variable enough
- 4 Match sequencing technology size
- 5 Represented in reference libraries

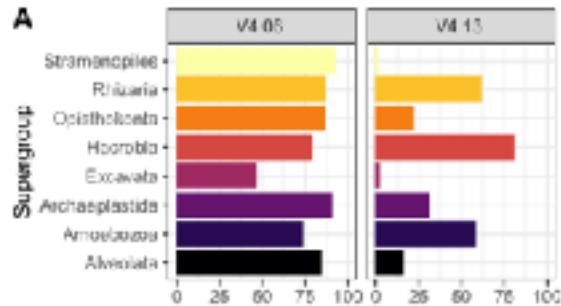
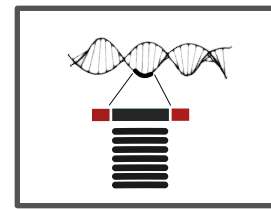
Match sequencing technology size



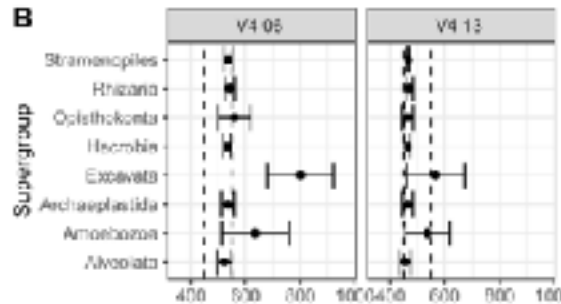
ADNr 18S : 1800bp
 ADNr 16S : 1500bp
 RbcL : 1400bp
 MatK : 1500bp
 COI : 650bp



Match sequencing technology size

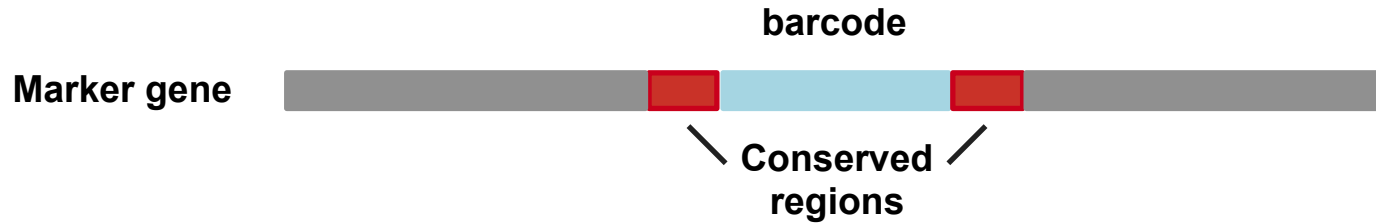
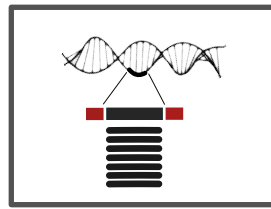


For the same primers, the amplicon size will vary between major supergroup of protists



A broader diversity of protist will be recovered with long-read sequencing

Barcode selection

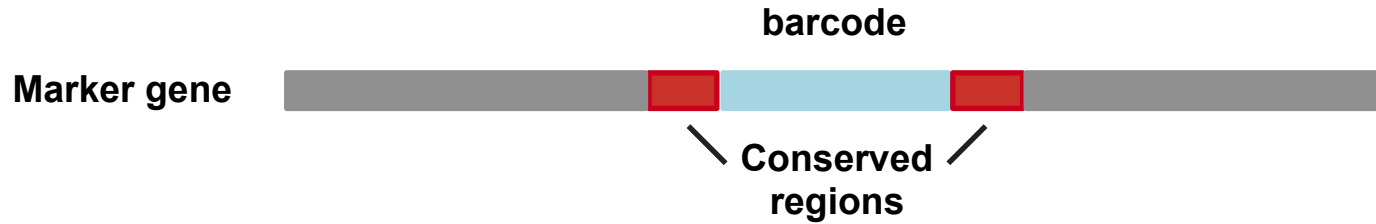
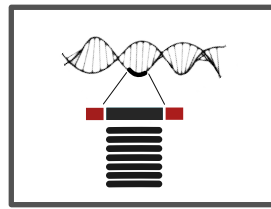


- 1 Universal
- 2 Conserved regions
- 3 Variable enough
- 4 Match sequencing technology size
- 5 Represented in reference libraries

The perfect barcode does not exist

Particularly if you want to target very diverse groups of organisms (e.g. all the protists, all the macroinvertebrates, all the phytoplankton)

Barcode selection

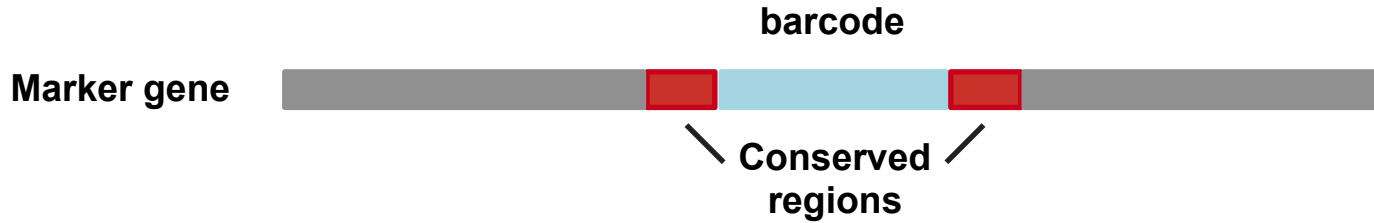
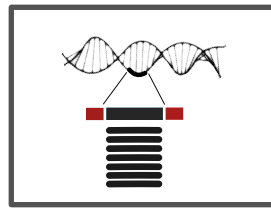


- 1 Universal
- 2 Conserved regions
- 3 Variable enough
- 4 Match sequencing technology size
- 5 Represented in reference libraries

The perfect barcode does not exist

Find the one that best suits to your research question

Barcode selection



- 1 Universal
- 2 Conserved regions
- 3 Variable enough
- 4 Match sequencing technology size

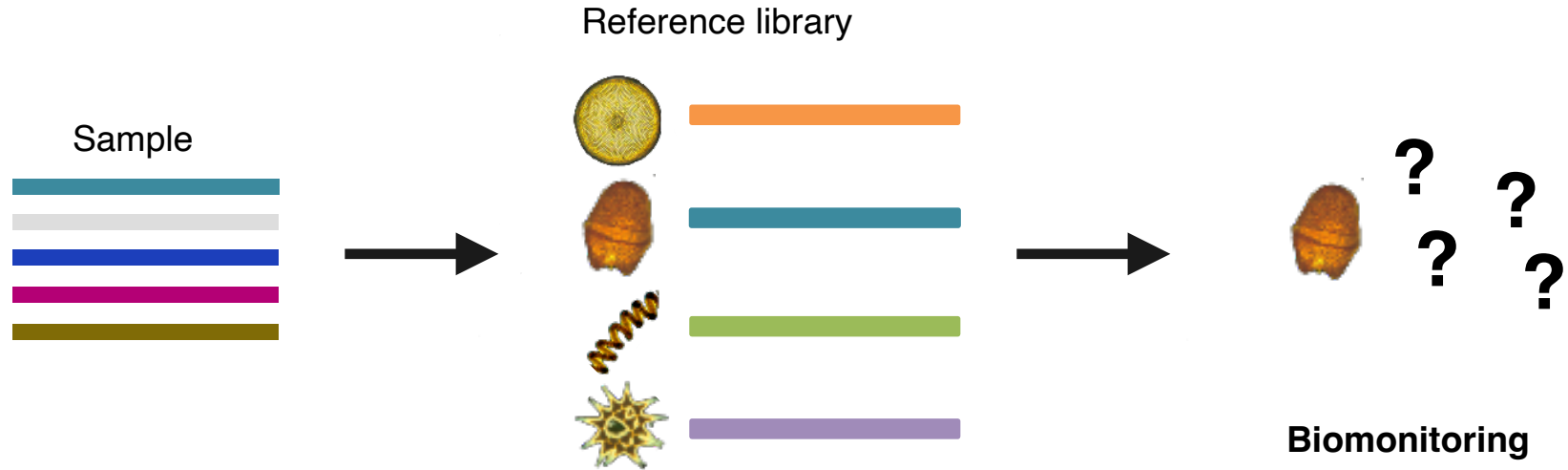
« *The power of DNA metabarcoding is directly proportional to data available in reference databases* »

- 5 Represented in reference libraries


Reference libraries



They are a prerequisite to identify your sequence and give it a taxonomy



Reference libraries

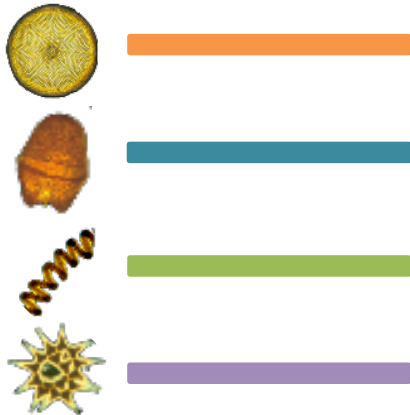
 ATCGCTTTGGACCT

 ATCGCTTTGGACCT

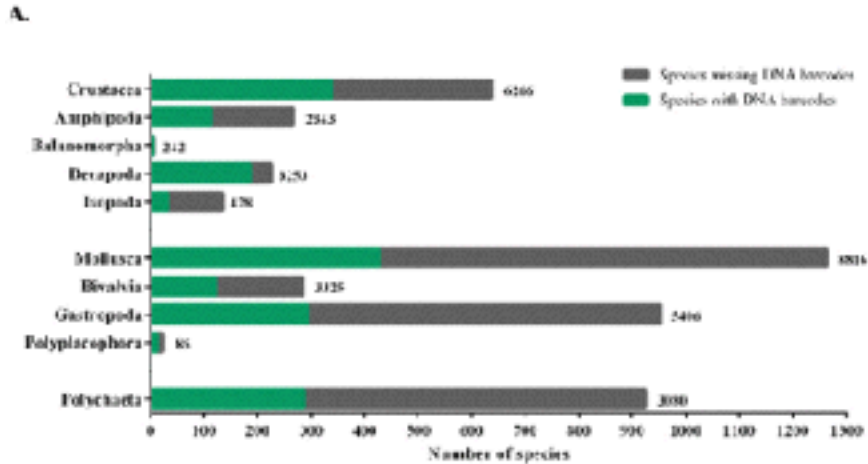
 ATCGCTTTGGACCT

They are a prerequisite to identify your sequence and give it a taxonomy

Reference library

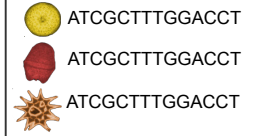


→ The most complete as possible



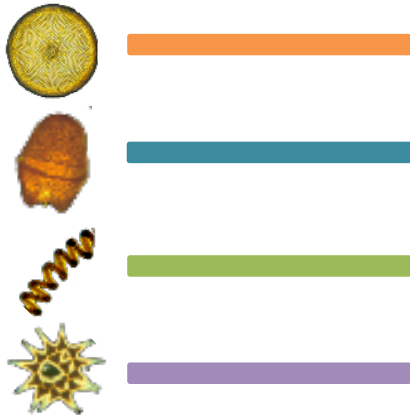
Leite et al., 2020

Reference libraries



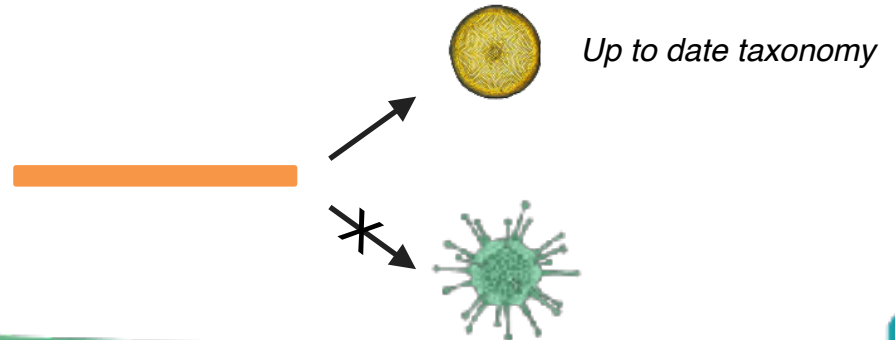
They are a prerequisite to identify your sequence and give it a taxonomy

Reference library



→ The most complete as possible

→ The most curated as possible





Generalists



SILVA SSU 138.1 update release

	SSU Parc	SSU Ref NR 99	LSU Parc	LSU Ref NR 99
Minimal length	300	1200/900	300	1900
Quality filtering	basic	strong	basic	strong
Guide Tree	no	yes	no	yes
Release date	27.08.20	27.08.20	27.08.20	27.08.20
Aligned rRNA sequences	9,469,124	510,508	1,312,534	95,286

For bacteria, archaea and eukaryotes
Both nuclear and plastid SSU and LSU

Specialists



Nuclear
18S

200,000
sequences

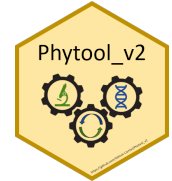
For protists

PhytoRef

Plastidal
16S

6486
sequences

For phytoplankton



Nuclear 18S
Plastidal
16S, 23S

18S : 16654
16S : 8479
23S : 1997


Generalists

BOLD
SYSTEMS



BOLD: The Barcode of Life Data System
(www.barcodinglife.org)

SUJEEVAN RATNASINGHAM and PAUL D. N. HEBERT

Around 1.3M of COI barcodes



Specialists

Diatbarcode	COIs
	
rbcl	COI
4000-5000 sequences	532 000 sequences
Diatoms	Insects

Reference libraries - non curated database

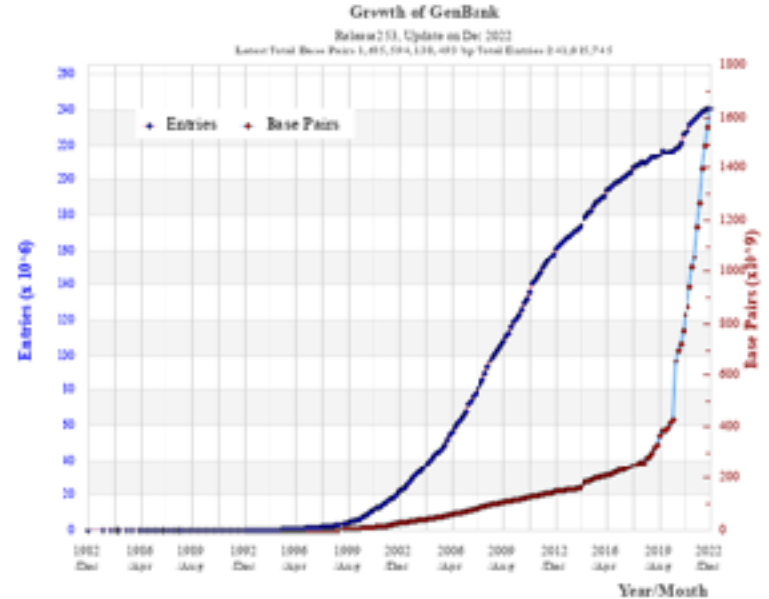
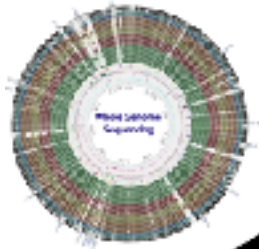


Your study



Big eDNA projects

WGS



For GenBank : 240 millions of sequences

High potential but **not curated** !

Reference libraries



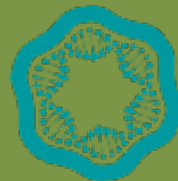
All sequences



Quality-checked



Curated taxonomy



PhytoRef

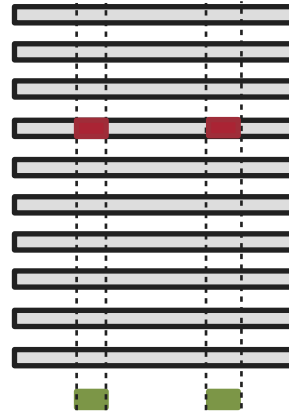
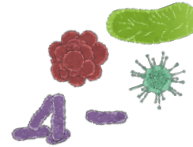
Reference libraries



But also they are essential to find the best design for the primers



Optimal design of the primers



Test for primer specificity

Reference libraries



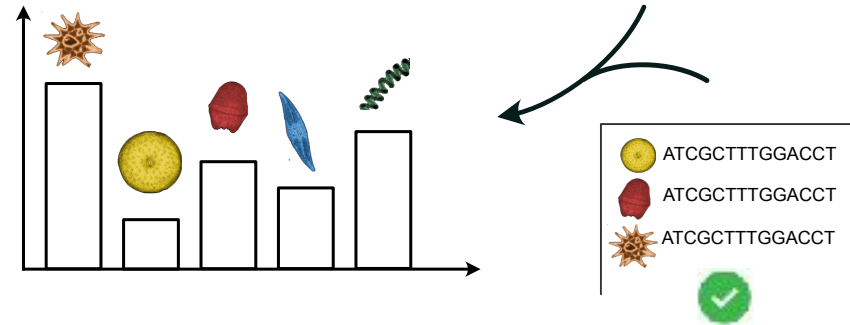
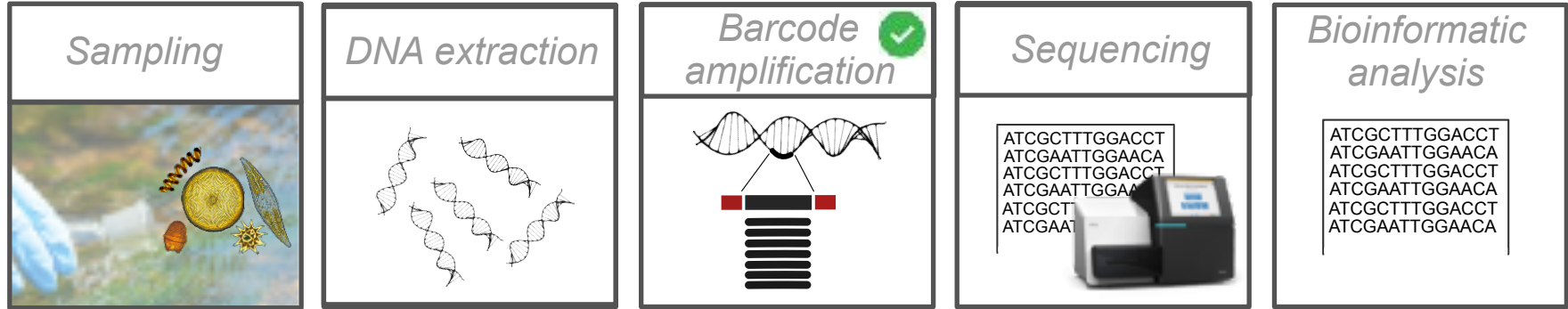
But also they are essential to find the best design for the primers

Optimal design of primers



st for primer specificity

Metabarcoding steps





Any questions?

