# Marine Ecological Genetics

## 12. DNA barcoding | Computer practical

- Extract DNA barcodes from Sanger reads

- Identify samples with the Barcode of Life Data System

- Evaluate genetic distances and id quality

Martin Helmkampf

Summer 2024

Exercises in Marine Ecological Genetics
12. DNA barcoding | Computer practical

Carl von Ossietzky
Universität
Oldenburg

# Access files for practical

Download files from GitHub:

https://github.com/mhelmkampf/meg24.git (Code | Download ZIP)

Alternatively, run git from terminal:

```
git clone https://github.com/mhelmkampf/meg24.git
```

Open R script called **Mar_Ecol_Gen_week2.R** in RStudio

Summer 2024

Exercises in Marine Ecological Genetics
12. DNA barcoding | Computer practical

Carl von Ossietzky
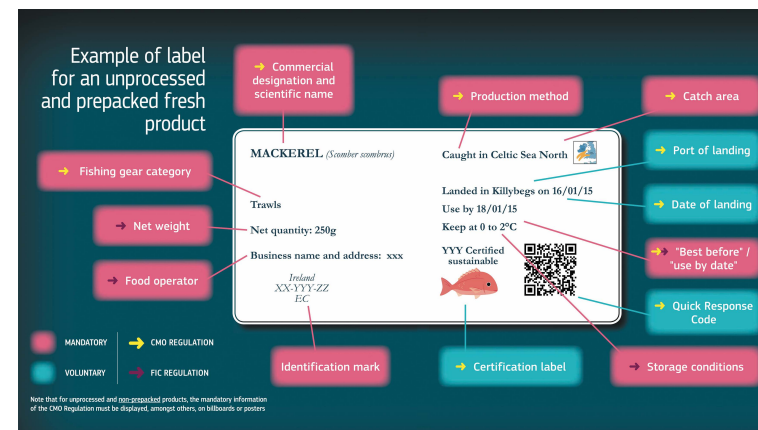Universität
Oldenburg

# #fischdetektive

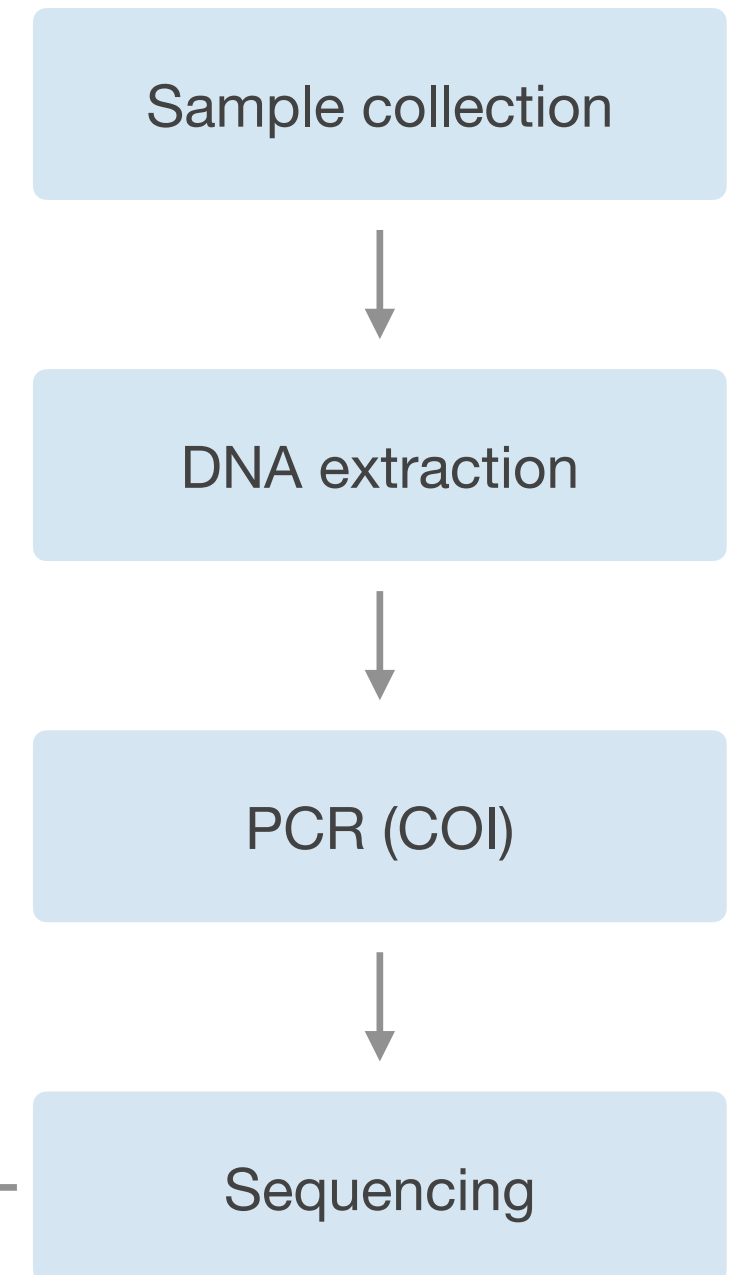Citizen science project at GEOMAR (2017) with over 700 participants (10–14 years)

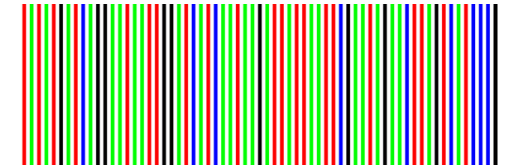Where does our seafood come from, and is it labeled correctly?

Thorsten Reusch, GEOMAR



Example of label for an unprocessed and prepacked fresh product

Sample collection

→ DNA extraction
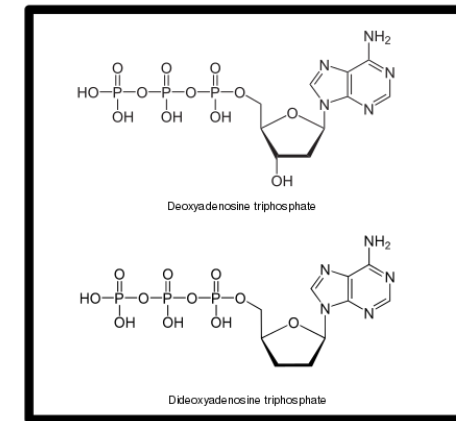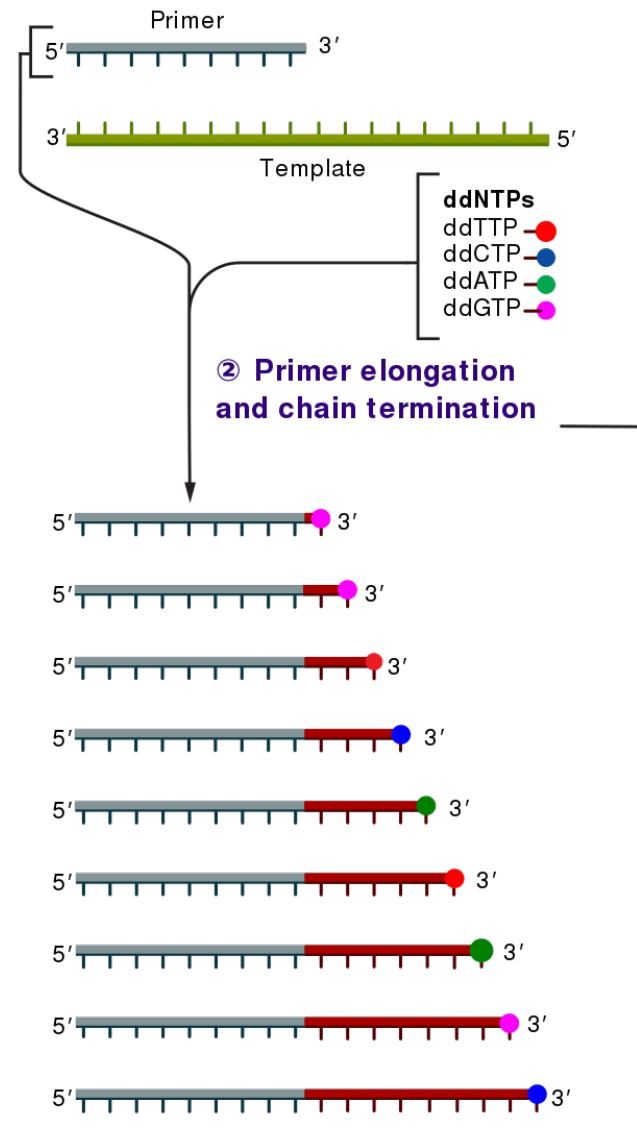
→ PCR (COI)

→ Sequencing → ID

# COI barcode

Approx. 650 bp in 5' region of cytochrome c oxidase subunit I (COI)

```
>MN604318.1 Oncorhynchus keta cytochrome c oxidase subunit I gene, complete cds; mitochondrial
GTGGCAATCACACGATGATTCTTCTCAACCAACCACAAAGACATTGGCACCCTCTATTTAGTATTTGGTGCCTGAGCCGGGATAGTAGGCACCGCCCTG
AGCCTACTAATTCGGGCAGAACTAAGCCAGCCAGGCGCTCTTCTAGGGGATGACCAGATCTACAATGTAATCGTTACAGCCCATGCCTTCGTTATAATT
TTCTTTATAGTCATACCAATTATAATCGGAGGCTTTGGAAACTGATTAATCCCCCTAATGATCGGGGCACCAGATATAGCATTCCCACGAATAAATAAC
ATAAGCTTCTGACTCCTACCTCCGTCCTTCCTCCTCCTCCTTTCTTCATCTGGAGTTGAAGCCGGCGCTGGTACCGGGTGGACAGTTTATCCCCCTCTA
GCCGGAAACCTTGCCCACGCAGGAGCATCTGTCGACTTAACCATCTTCTCCCTCCATTTAGCTGGAATCTCCTCAATTTTGGGGGCCATTAATTTTATT
ACGACCATTATCAACATAAAACCCCCAGCTATTTCTCAGTACCAAACCCCGCTTTTTGTCTGAGCTGTACTAATCACTGCTGTACTTCTACTATTATCA
CTCCCCGTTCTGGCAGCAGGTATTACTATGTTGCTCACAGATCGAAATTTAAACACCACTTTCTTTGACCCGGCGGGTGGCGGAGATCCAATTTTATAC
CAACACCTCTTTTGATTCTTCGGTCACCCAGAGGTCTATATTCTGATCCTCCCAGGCTTTGGTATAATTTCACATATCGTTGCATATTACTCTGGTAAG
AAAGAACCTTTCGGGTACATAGGAATAGTGTGAGCTATAATAGCCATCGGCTTGTTAGGATTTATCGTTTGAGCCCACCACATATTTACTGTCGGGATG
GACGTGGACACTCGTGCCTACTTTACATCTGCCACCATAATTATCGCTATCCCCACAGGAGTAAAAGTATTTAGCTGACTAGCTACACTGCACGGAGGC
TCGATCAAATGAGAGACACCACTTCTCTGAGCCCTAGGATTTATCTTCCTATTTACAGTGGGCGGATTAACGGGCATCGTCCTTGCTAACTCCTCATTA
GACATTGTTTTACATGACACTTATTACGTAGTCGCCCATTTCCACTACGTACTCTCAATAGGAGCTGTATTTGCCATTATGGGCGCTTTCGTACACTGA
TTCCCCCTATTCACAGGGTACACCCTTCACAGCACATGAACCAAAATCCATTTTGGAATTATATTTATCGGTGTAAATTTAACCTTTTTCCCACAGCAT
TTCCTAGGCCTCGCAGGGATACCACGACGGTACTCTGACTACCCGGACGCCTACACGCTATGAAACACTGTATCCTCAATCGGATCCCTTGTCTCCTTA
GTAGCTGTAATTATGTTCCTATTTATTCTTTGAGAGGCTTTTGCTGCCAAACGAGAAGTAGCATCAATCGAAATAACTTCAACAAACGTAGAATGACTA
CACGGATGCCCCCCACCCTACCACACATTCGAGGAACCAGCATTTGTCCAAGTACGAACGTACTAA
```
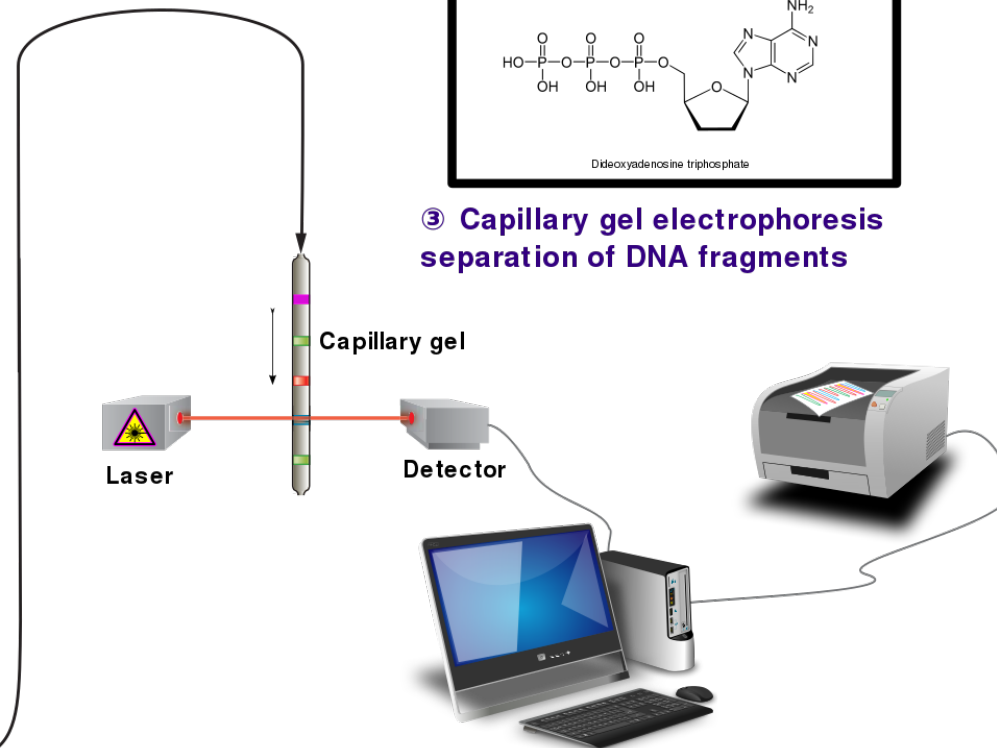
Exercises in Marine Ecological Genetics

12. DNA barcoding | Computer practical

Carl von Ossietzky
Universität
Oldenburg

# Sanger sequencing



① Reaction mixture
- Primer and DNA template
- DNA polymerase
- ddNTPs with flourochromes
- dNTPs (dATP, dCTP, dGTP, and dTTP)

Primer
5'                                3'

3'                                5'
Template

**ddNTPs**
ddTTP ●
ddCTP ●
ddATP ●
ddGTP ●

② Primer elongation and chain termination

③ Capillary gel electrophoresis separation of DNA fragments

Deoxyadenosine triphosphate

Dideoxyadenosine triphosphate

Capillary gel

Laser

Detector

④ Laser detection of flourochromes and computational sequence analysis

Chromatograph

Estevezj, CC BY-SA 3.0

Exercises in Marine Ecological Genetics

12. DNA barcoding | Computer practical

Carl von Ossietzky
Universität
Oldenburg

# Sanger read processing

Basecalling

Trace file (F + R)

Trimming

Fasta file (F+R)

Reverse complement R
Align F and R

R ‖‖‖‖‖‖‖‖‖‖‖‖ F

Alignment

Create consensus sequence

Consensus

Exercises in Marine Ecological Genetics
12. DNA barcoding | Computer practical

Carl von Ossietzky
Universität
Oldenburg

# Sequence alignment

```
G A T G T T C G A A
G A T C – – – G A A
G A C C – T C G – T
```

Arranges nucleotide or amino acid sequences

so that the number of mismatches and gaps are minimized

- Multiple sequence alignments can be constructed progressively from pairwise alignments

- Computationally complex, often requires heuristic solutions

- Key to identify evolutionary relationships between sequences (e.g. homology)

Summer 2024

Exercises in Marine Ecological Genetics
12. DNA barcoding | Computer practical

Carl von Ossietzky
Universität
Oldenburg

# Genetic distances

```
AAGCCAGCCAGGCGCTCTTCTAGGGGATGACCAGATCTACAATGTAATCG    # 50 positions total

AAGTCAACCTGGTGCACTTCTTGGTGATGATCAAATTTATAATGTGATCG

***.**.** **.** ***** ** ** **.**.**.**.*****.****    # 13 differences
```

Uncorrected distance

$p$ = 13 / 50 = 0.26

K2P distance (Kimura 1980)

$K$ = 0.33

$$K = -\frac{1}{2}\ln\left((1 - 2p - q)\sqrt{1 - 2q}\right)$$

$p$: proportion of transitions (A↔G, C↔T)

8 / 50 = 0.16

$q$: proportion of transversions

5 / 50 = 0.1

Exercises in Marine Ecological Genetics
12. DNA barcoding | Computer practical

Carl von Ossietzky
Universität
Oldenburg

# BOLD

Exercises in Marine Ecological Genetics
12. DNA barcoding | Computer practical

Carl von Ossietzky
Universität
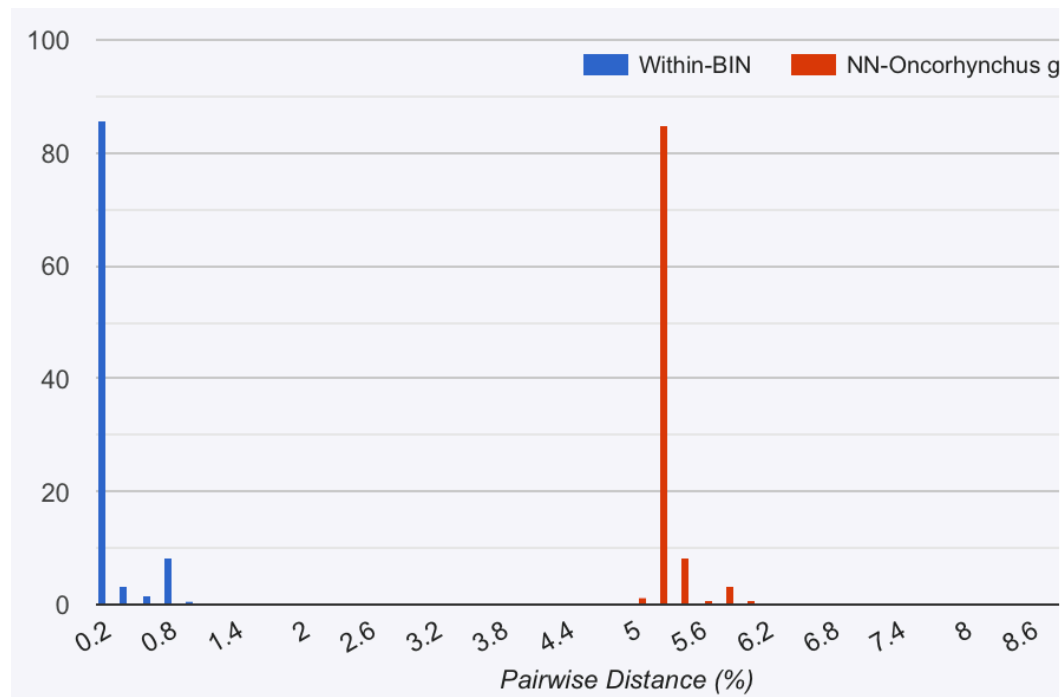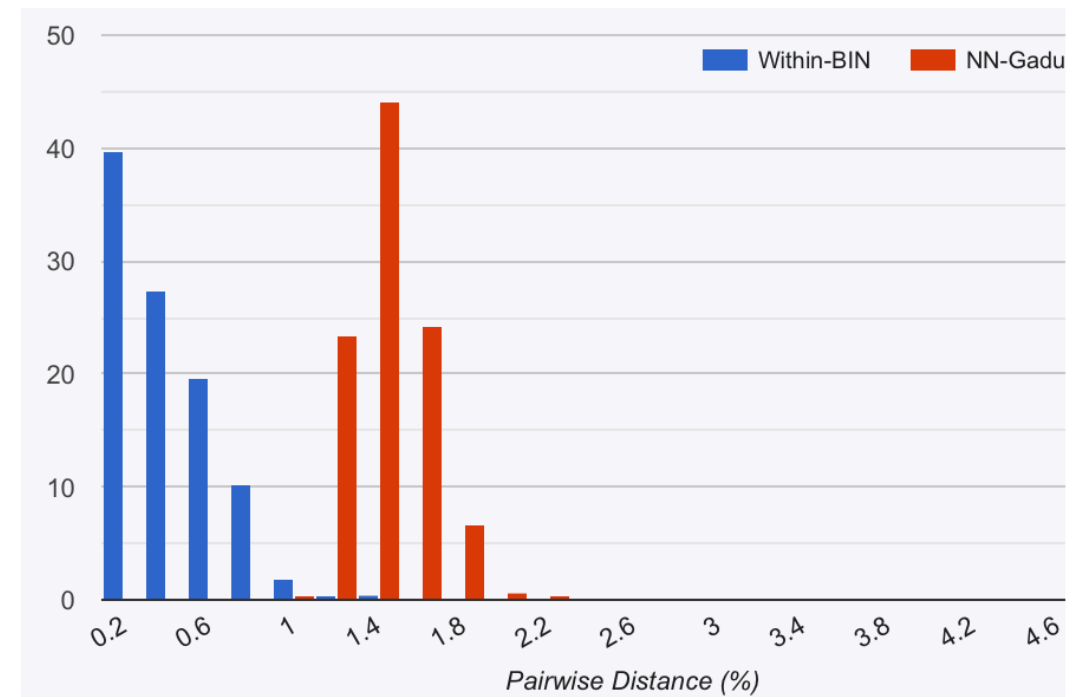Oldenburg

# Barcode gap

Comparing genetic distances within BIN and to nearest neighbor (NN) BIN



Large barcode gap

*Oncorhynchus keta*

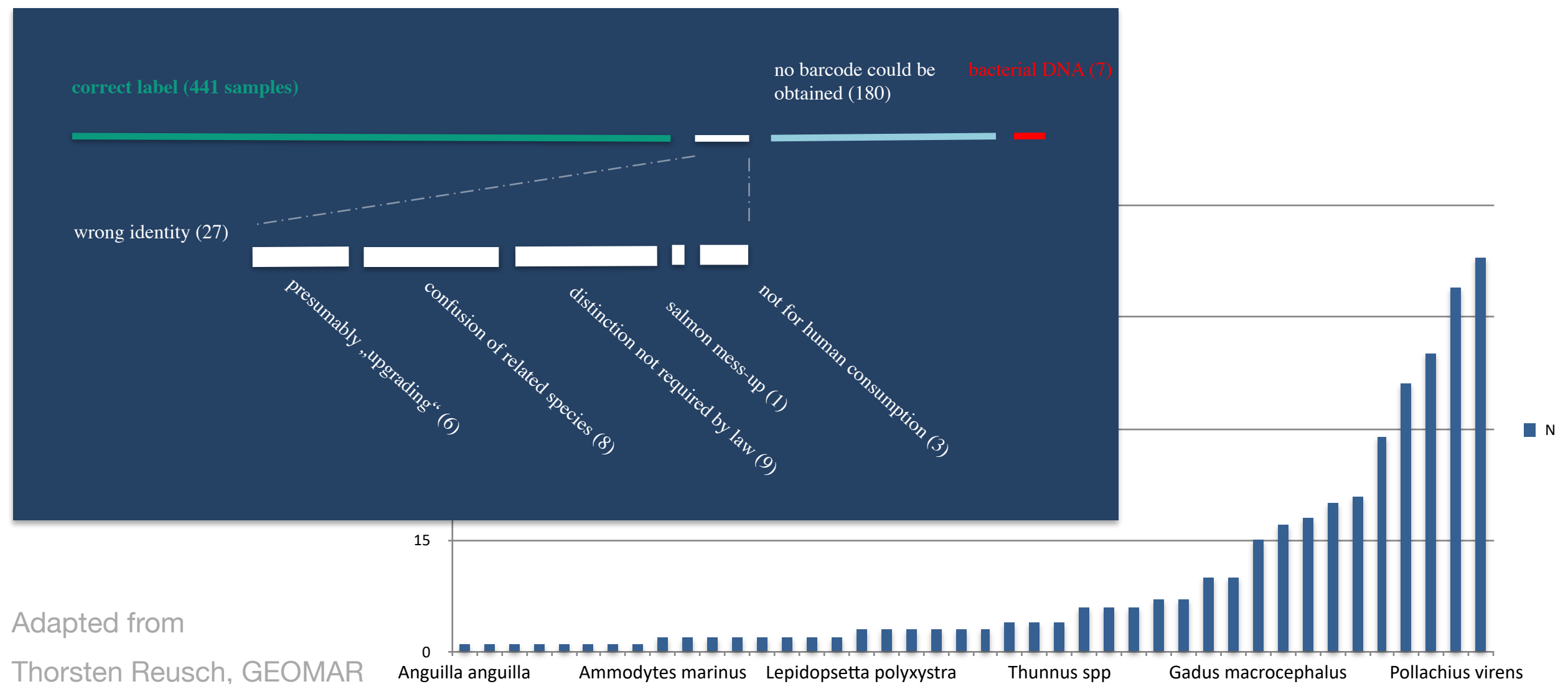Small barcode gap

*Gadus chalcogrammus*

Summer 2024

Exercises in Marine Ecological Genetics
12. DNA barcoding | Computer practical

Carl von Ossietzky
Universität
Oldenburg

# #fischdetektive results

Mislabeling seems to be only a moderate problem in Germany in frozen and fresh fish

(but may be higher in Sushi-grade fish and processed fish products)



correct label (441 samples)

no barcode could be obtained (180)

bacterial DNA (7)

wrong identity (27)

presumably „upgrading" (6)

confusion of related species (8)

distinction not required by law (9)

salmon mess-up (1)

not for human consumption (3)

15

0

Anguilla anguilla   Ammodytes marinus   Lepidopsetta polyxystra   Thunnus spp   Gadus macrocephalus   Pollachius virens

N

Summer 2024

Exercises in Marine Ecological Genetics

12. DNA barcoding | Computer practical

Carl von Ossietzky
Universität
Oldenburg

# Portable 3rd gen sequencing

Oxford Nanotechnologies MinION



96-well plate

Indexing during PCR
Library preparation



whatech.com



NANOPORE SEQUENCING

At the heart of the MinION device, an enzyme unwinds DNA, feeding one strand through a protein pore. The unique shape of each DNA base causes a characteristic disruption in electrical current, providing a readout of the underlying sequence.

DNA double helix

DNA base

Unwinding enzyme

Protein pore

Membrane

Ion

Current

Current

Sequence    A  A  C  T  C  G  T

blogs.nature

Summer 2024

Exercises in Marine Ecological Genetics
12. DNA barcoding | Computer practical

Carl von Ossietzky
Universität
Oldenburg

# Take-home messages

- Extract DNA barcodes from Sanger reads

- Identify samples with the Barcode of Life Data System

- Evaluate genetic distances and id quality

Summer 2024

Exercises in Marine Ecological Genetics
12. DNA barcoding | Computer practical

Carl von Ossietzky
Universität
Oldenburg