# Exercises in Marine Ecological Genetics

## 09. DNA barcoding

- Extract barcodes from Sanger reads

- Match sequence to BOLD database

- Evaluate id quality using genetic distances

Martin Helmkampf



Foto: dpa

https://github.com/mhelmkampf/meg25

Carl von Ossietzky
Universität
Oldenburg

# #fischdetektive

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

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





Thorsten Reusch, GEOMAR
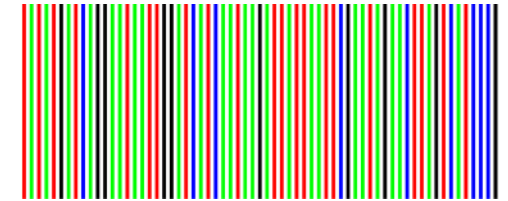
Sample collection

↓

DNA extraction

↓

PCR (COI)

↓

Matching  ←  Sequencing

Summer 2025

Exercises in Marine Ecological Genetics
09. DNA barcoding

Carl von Ossietzky
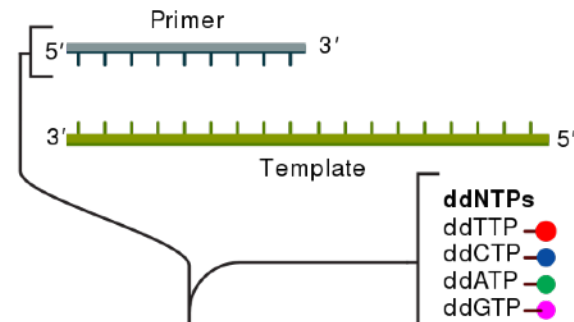Universität
Oldenburg

# 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
```
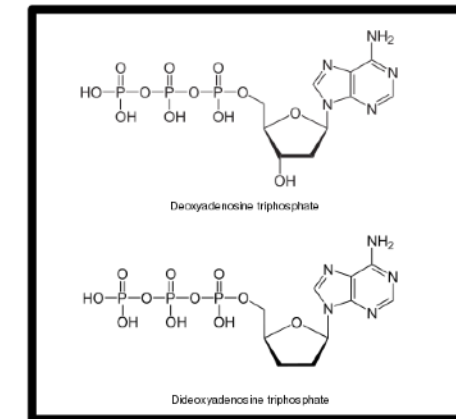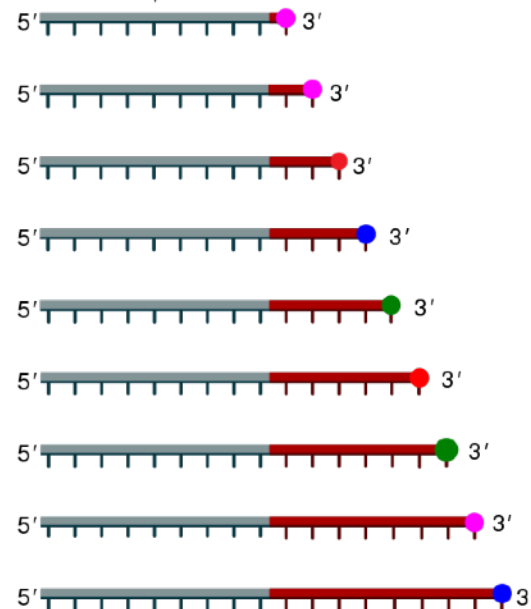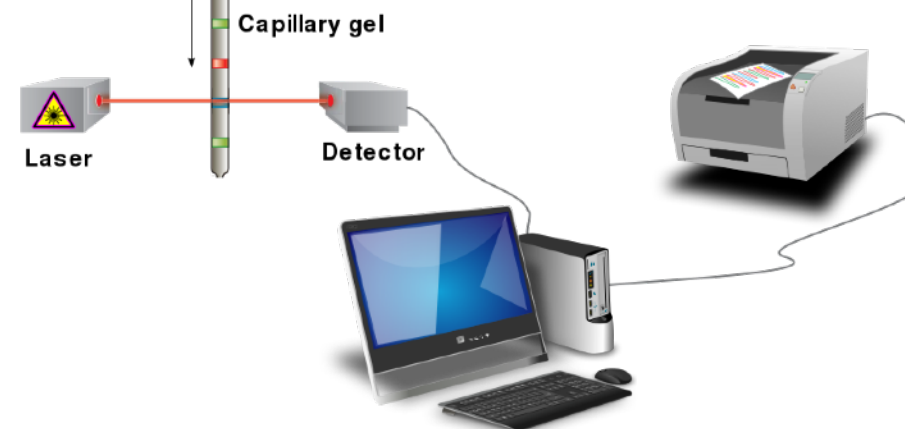
Summer 2025

Exercises in Marine Ecological Genetics
09. DNA barcoding

Carl von Ossietzky
Universität
Oldenburg

# Sanger sequencing



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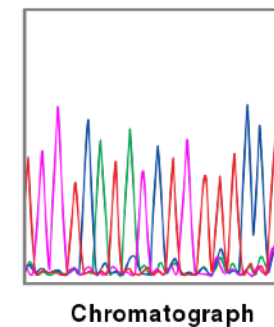‣ **ddNTPs with flourochromes**‣ **dNTPs (dATP, dCTP, dGTP, and dTTP)**

Primer
5′ —————— 3′

3′ —————————— 5′
Template

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

② **Primer elongation and chain termination**

Decoxyadenosine triphosphate

Dideoxyadenosine triphosphate

③ **Capillary gel electrophoresis separation of DNA fragments**

Capillary gel

Laser        Detector

④ **Laser detection of flourochromes and computational sequence analysis**

Chromatograph

Exercises in Marine Ecological Genetics
09. DNA barcoding

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

**Alignment**

Create consensus sequence

**Consensus**

R | | | | | | | | | | | | F

Exercises in Marine Ecological Genetics
09. DNA barcoding

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 are minimized

- Accomplished by introducing gaps (–), which represent insertions or deletions (**indels**) and account for sequence length differences due to mutations over time

- Computationally complex, often requires heuristic solutions

- **Reveals evolutionary or functional relationships** between sequences (e.g. homology)

- Key to variant calling, sequence assembly, species identification and phylogenetics

Summer 2025

Exercises in Marine Ecological Genetics
09. DNA barcoding

Carl von Ossietzky
Universität
Oldenburg

# Barcode Index Number (BIN)

- **Clusters of similar COI barcodes**

- Clustering is based on genetic similarity, independent of formal taxonomy

- Represent operational taxonomic units (OTUs) that often correspond to species

- Enable species identification and biodiversity assessments

Exercises in Marine Ecological Genetics
09. DNA barcoding

Carl von Ossietzky
Universität
Oldenburg

# Genetic distance

Uncorrected or *p*-distance:

Proportion of nucleotides at which two sequences differ

```
AAGCCAGCCAGGCGCTCTTCTAGGGGATGACCAGATCTACAATGTAATCG    # 50 positions total

AAGTCAACCTGGTGCACTTCTTGGTGATGATCAAATTTATAATGTGATCG

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

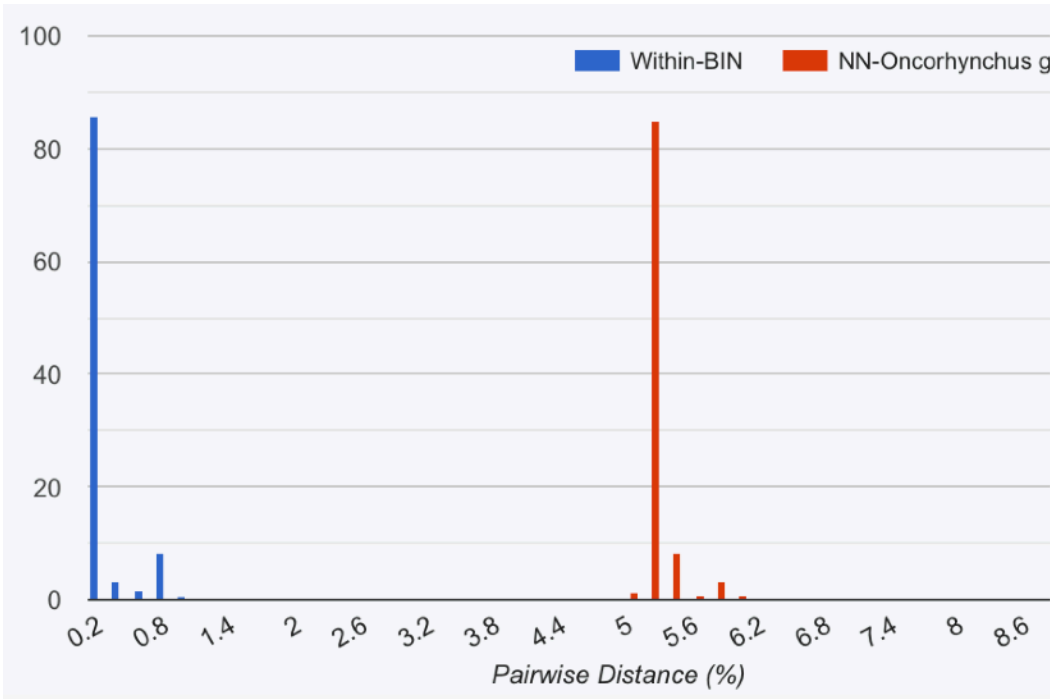*p* = 13 / 50 = 0.26

Does not correct for multiple substitutions (repeated mutations of the same site)

or transition / transversion bias (> K2P distance, Kimura 1980)

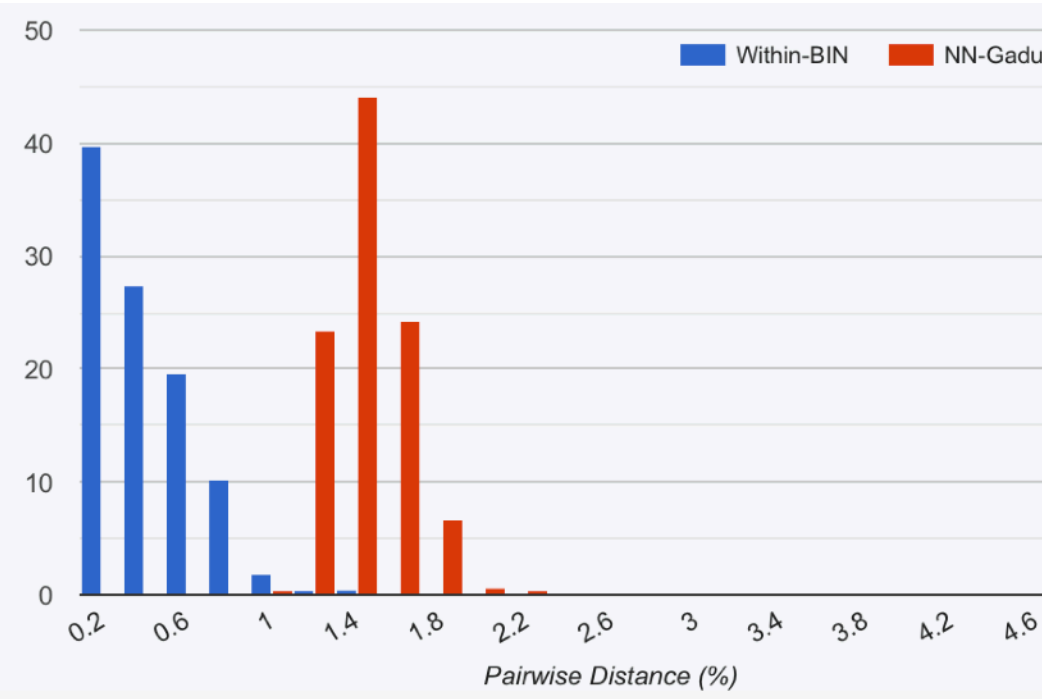Carl von Ossietzky
Universität
Oldenburg

# Barcode gap

Comparing genetic distances within BIN

and to nearest neighbor (NN) BIN

| Comparison | Typical p-distance (COI) |
|---|---|
| Within species | < 2 % |
| Between species | > 2–3 % |
| Between genera | 10–20 % |



Large barcode gap

*Oncorhynchus keta*



Small barcode gap

*Gadus chalcogrammus*

Summer 2025

Exercises in Marine Ecological Genetics
09. DNA barcoding
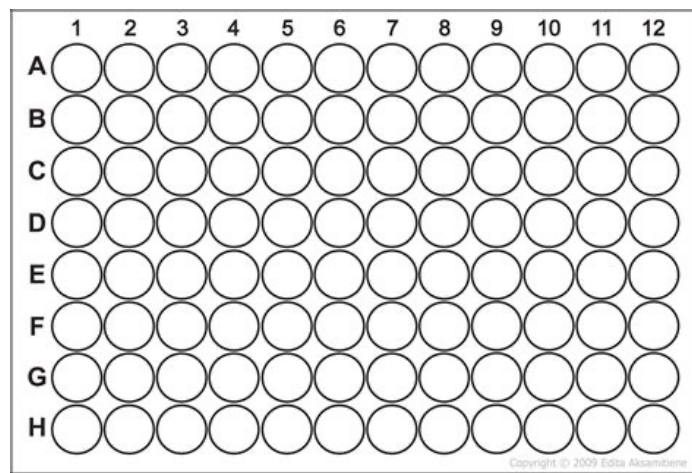
Carl von Ossietzky
Universität
Oldenburg

# #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)



Adapted from

Thorsten Reusch, GEOMAR

Exercises in Marine Ecological Genetics

09. DNA barcoding

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.

blogs.nature

Summer 2025

Exercises in Marine Ecological Genetics
09. DNA barcoding

Carl von Ossietzky
Universität
Oldenburg

# DNA barcoding

```
seqkit subseq –r start:end input.fas > output.fas

seqkit seq –r –p input.fas > output.fas

merger –asequence forward.fas –bsequence reverse.fas … –outseq consensus.fas
```

- Accurate barcodes like COI can be extracted from high-quality Sanger reads after **careful processing**, including trimming and consensus formation

- **Matching to a reference database** like BOLD allows to assign samples to species, but results depend on sequence quality and database completeness

- **Identification reliability** may be assessed using sequence similarity and genetic distance to nearest-neighbor ("barcode gap")

Summer 2025

Exercises in Marine Ecological Genetics
09. DNA barcoding

Carl von Ossietzky
Universität
Oldenburg

# Evaluation — please participate!



https://elearning.uni-oldenburg.de/plugins.php/unizensusplugin/show?cid=3660d16e8eb3daf479389cf8233c12fb

Summer 2025

Exercises in Marine Ecological Genetics
09. DNA barcoding

Carl von Ossietzky
Universität
Oldenburg