
Rapid species level identification of fish eggs by proteome fingerprinting using MALDI-TOF MS

Rossel Sven ^{1,*}, Barco Andrea ², Kloppmann Matthias ³, Martínez Arbizu Pedro ¹, Huwer Bastian ⁴, Knebelberger Thomas ²

¹ enckenberg am Meer, German Centre for Marine Biodiversity Research (DZMB), Südstrand 44, 26382 Wilhelmshaven, Germany

² biome-ID, Südstrand 44, 26382 Wilhelmshaven, Germany

³ Thünen Institut für Seefischerei, Herwigstraße 31, 27572, Bremerhaven, Germany

⁴ Technical University of Denmark, National Institute of Aquatic Resources, Kemitorvet, Bygning 202, 2800 Kgs. Lyngby, Denmark

* Corresponding author : Sven Rossel, email address : sven.rossel@senckenberg.de

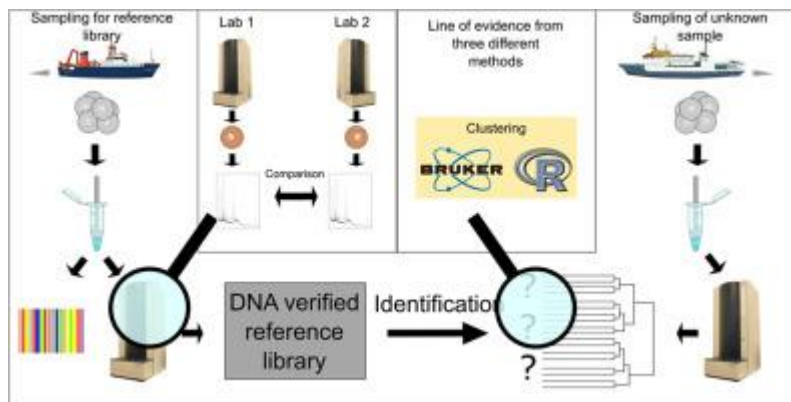
Abstract :

Quantifying spawning biomass of commercially relevant fish species is important to generate fishing quotas. This will mostly rely on the annual or daily production of fish eggs. However, these have to be identified precisely to species level to obtain a reliable estimate of offspring production of the different species. Because morphological identification can be very difficult, recent developments are heading towards application of molecular tools. Methods such as COI barcoding have long handling times and cause high costs for single specimen identifications. In order to test MALDI-TOF MS, a rapid and cost-effective alternative for species identification, we identified fish eggs using COI barcoding and used the same specimens to set up a MALDI-TOF MS reference library. This library, constructed from two different MALDI-TOF MS instruments, was then used to identify unknown eggs from a different sampling occasion. By using a line of evidence from hierarchical clustering and different supervised identification approaches we obtained concordant species identifications for 97.5% of the unknown fish eggs, proving MALDI-TOF MS a good tool for rapid species level identification of fish eggs. At the same time we point out the necessity of adjusting identification scores of supervised methods for identification to optimize identification success.

Significance :

Fish products are commercially highly important and many societies rely on them as a major food resource. Over many decades stocks of various relevant fish species have been reduced due to unregulated overfishing. Nowadays, to avoid overfishing and threatening of important fish species, fish stocks are regularly monitored. One component of this monitoring is the monitoring of spawning stock sizes. Whereas this is highly dependent on correct species identification of fish eggs, morphological identification is difficult because of lack of morphological features.

Graphical abstract



Highlights

► Application of proteome fingerprinting for fish stock monitoring. ► 97.5% identification success. ► Adjustment of identification scores improves identification success. ► Line of evidence from different methods improves identification confidence. ► Automatic mass spectra quality control.

Keywords : Species identification, Fishing quota, Inter laboratory, Random forest, Bruker Biotyper, Fishery, COI

Introduction

The correct identification of fish eggs is an important prerequisite in research of reproduction ecology of fish species [1–4] but also for the estimation of spawning stock size of economically important fish species for assessment purposes [5–9]. However, identification of fish eggs is notoriously difficult, particularly if they are of an early developmental stage. These early stages are most often the only stages that are used to calculate egg production for spawning stock biomass estimation [1,5,8]. Apart from species of Callionymids, Beloniformes, Macrourids or Mauroliius, which have a characteristically shaped chorion, the early eggs of the vast majority of marine fish species come with no other morphological distinguishing feature than the egg diameter and presence as well as size and number of oil droplets in the yolk [10,11]. With a typical range of 0.6 – 2.0 mm of egg diameter and 0.1 – 0.4 mm for the oil droplet [10,11], there is a considerable overlap in those measures among the several marine fish species, making it almost impossible to determine the species of newly spawned eggs [12]. This is e.g. the case in the egg survey for the winter spawning fish in the North Sea [13], where cod eggs have to be separated from all other similar sized gadoid eggs or during the mackerel and horse mackerel egg survey [14], where mackerel (*Scomber scombrus* Linnaeus, 1758) eggs have to be distinguished from the similar sized hake (*Merluccius merluccius* (Linnaeus, 1758)) and ling (*Molva molva* (Linnaeus, 1758)) eggs [15]. Due to these challenges in morphological identification, a number of studies have recently focused on fish egg species identification using molecular methods such as COI barcoding [3,4,16–9] with overall good success.

Another method for rapid and reliable species identification is Matrix-Assisted Laser Desorption/Ionization Time-of-Flight Mass Spectrometry (MALDI-TOF MS). This method relies on a so-called proteome fingerprint to distinguish between species [20]. Currently, this method is routinely applied in microorganism identification such as detection of bacterial pathogens [21–23], viruses [24,25] or fungi [26]. Aside from microbiology, nowadays MALDI-TOF MS is also applied to fight food fraud [27], e.g. to detect mislabeling of sea food species [28–35], identify meat origin [36] or inspect milk adulteration [37]. Moreover, MALDI-TOF MS was tested in numerous studies to identify important disease vectors such as mosquitos, ticks and phlebotomine sand flies [38–41]. But it was also successfully applied in ecological studies [42–44] with high species identification accuracy based on reference libraries. However, considering the need for accurate species identification in fisheries science described above, it is surprising that the method has so far not been applied to fish eggs.

The main advantages of this method opposing to DNA barcoding are the reduced costs [43] but also the short sample handling times while retaining high identification success. These advantages are particularly useful for a potential application in fish stock assessment, where relatively large numbers of eggs need to be identified and results usually need to be made available for assessment working groups with rather short deadlines. Therefore, the present study intends to test the suitability of MALDI-TOF MS for the identification of fish eggs, using field samples collected in the North Sea that contain a mixture of eggs of different marine fish species. At the same time we aim at comparing different identification strategies and recommending adjustments for identification thresholds of these different methods.

Material and Methods

Sampling:

Fish eggs for library construction were sampled with RV Walther Herwig III at different stations in the North Sea during WH413 between 22 January and 23 February 2018 (Fig. 1). Fish eggs for identification using MALDI-TOF MS were sampled with RV Dana on cruise number Dana/02/2018 between 01 and 19 February 2018. Eggs were sampled using the MIKey M net attachment [13] to the MIK net, which is deployed during the first quarter QBTS each year to catch large herring larvae [45]. The MIKey M net consists of a ring opening (20 cm diameter) and a 1.75 m long black 335 µm net. The MIKey M net, which is attached on the outside of the larger (2 m diameter) MIK net, is designed to catch small plankters including fish eggs concomitant with the MIK sampling down to a maximum depth of 100 m or 5 m above the seabed. The resulting catch is collected in a small codend, which can be detached and emptied upon retrieval of the net. For a detailed description of the MIKey M net sampling see ICES (2018) [13].

On RV Walther Herwig III the catch was immediately sorted for fish eggs in a tray placed on a bed of crushed ice in order to prevent quick deterioration of the eggs in the warm ship lab environment. The size of the eggs was measured, their developmental stage determined and then placed one by one into microcentrifuge tubes and stored in 96% undenatured ethanol at -18°C until further processing. On RV Dana, it was not possible to sort the catch on a bed of crushed ice. Instead, the sample was kept cold in a refrigerator and only small portions of the sample were sorted at a time under a stereo microscope. The eggs were sorted by developmental stage into small glass bowls, which were filled with chilled sea water and kept cold on blocks of ice. Each egg was then placed individually in a drop of sea water on a microscope slide with an engraved scale bar and photographed for later size measurements with an image analysis system. After the picture was taken, each egg was placed individually into an Eppendorf vial with 96% undenatured ethanol and stored at -18°C after all samples at a station were processed. Later on in the laboratory at land, eggs were transferred one by

one into microcentrifuge tubes for further processing. Storage time until processing did not exceed four months.

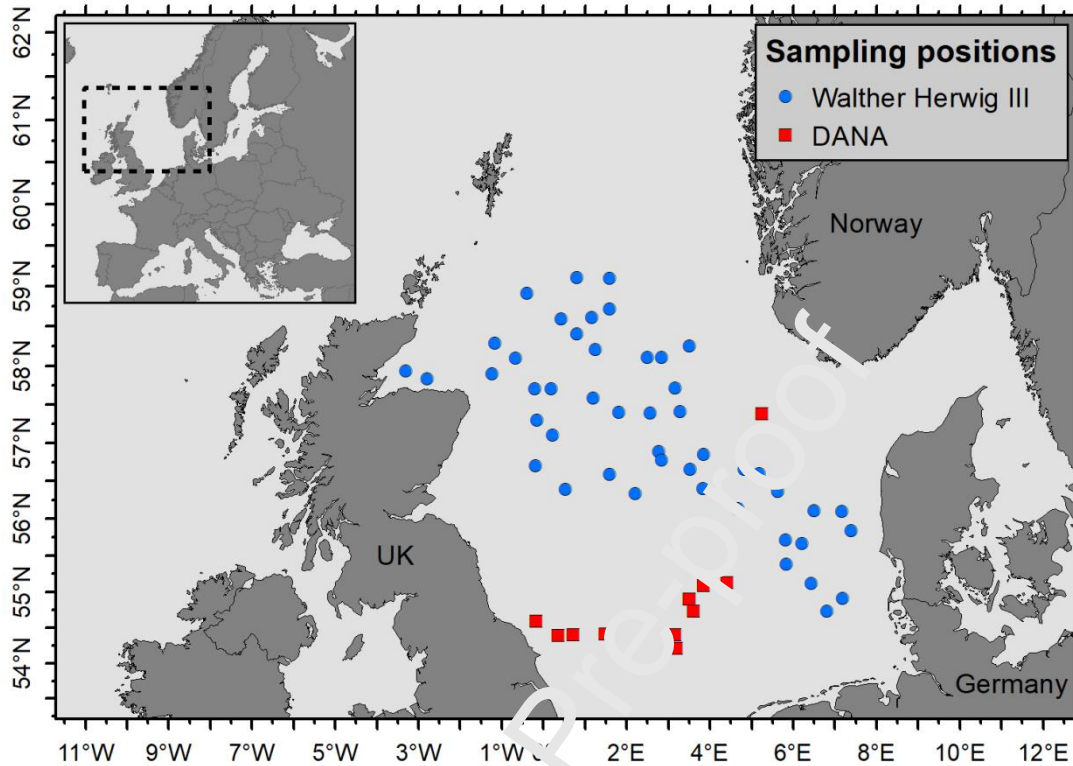


Fig. 1: Map of the different sampling sites across the North Sea with research vessels Walther Herwig III (blue circles) and Dana (red rectangles).

Sample processing

To allow preparations for COI barcoding and MALDI-TOF MS from the same individual egg, all individual eggs were transferred into 40 μ l molecular grade water and crushed using a microcentrifuge tube pestle (Fig. 2). Subsequently samples were vortexed and centrifuged for some seconds at 2000 g.

COI barcoding

Of the mixture, 35 μ l were used for DNA extraction using 200 μ l InstaGene matrix (Bio-Rad Laboratories, Munich, Germany). DNA extraction was only carried out for reference library eggs from WH413 (Fig. 2). For all other specimens, these 35 μ l were stored at -20 °C for later analyses. Genomic DNA was isolated at 95°C for 30 minutes (min) on a thermoblock. For amplification of the cytochrome-*c* oxidase I gene fragment, 5 μ l DNA were added to a mixture containing 18 μ l molecular grade water and 0.5 μ l of four different primers with a concentration of 10 pmol/ μ l each. A fish specific primer mix with M13 sequencing tails containing FR1d –t1 [46], VF2 – t1, FishF2 – t1 and

FishR2 – t1 [47] as described by Ivanova et al. (2007) [46] was used with illustra™ PuReTaq™ Ready-to-Go™ PCR beads. Initial denaturation was carried out at 95°C for 5 min. Following to this, 40 cycles including denaturation at 95°C for 1 min, annealing at 47°C for 1 min and elongation for 1 min 20 seconds were performed. Final elongation was done for 10 min at 72°C. Of the amplified PCR products, 2 µl were verified for size conformity by electrophoresis in a 1-% agarose gel stained with GelRED™ using commercial DNA size standards. Ten microliters of each PCR product were purified with a 2.5 µL mix containing exonuclease I (20 U/µL) and alkaline phosphatase (1 U/µL) using an incubation of 15 min at 37°C and 20 min at 75°C. Purified PCR products were sequenced unidirectional at a contract sequencing facility (Macrogen Europe, Amsterdam, Netherlands) using an ABI 3730xl DNA Sequencer and M13 universal primers. Sequencing results were quality controlled and resulting DNA sequences blasted to receive species identifications for reference eggs.

Data was aligned in SeaView [48] using muscle [49] algorithm and by eye control. A Neighbor-Joining tree was constructed using [50] based on Kimura-2 Parameter [51] distances.

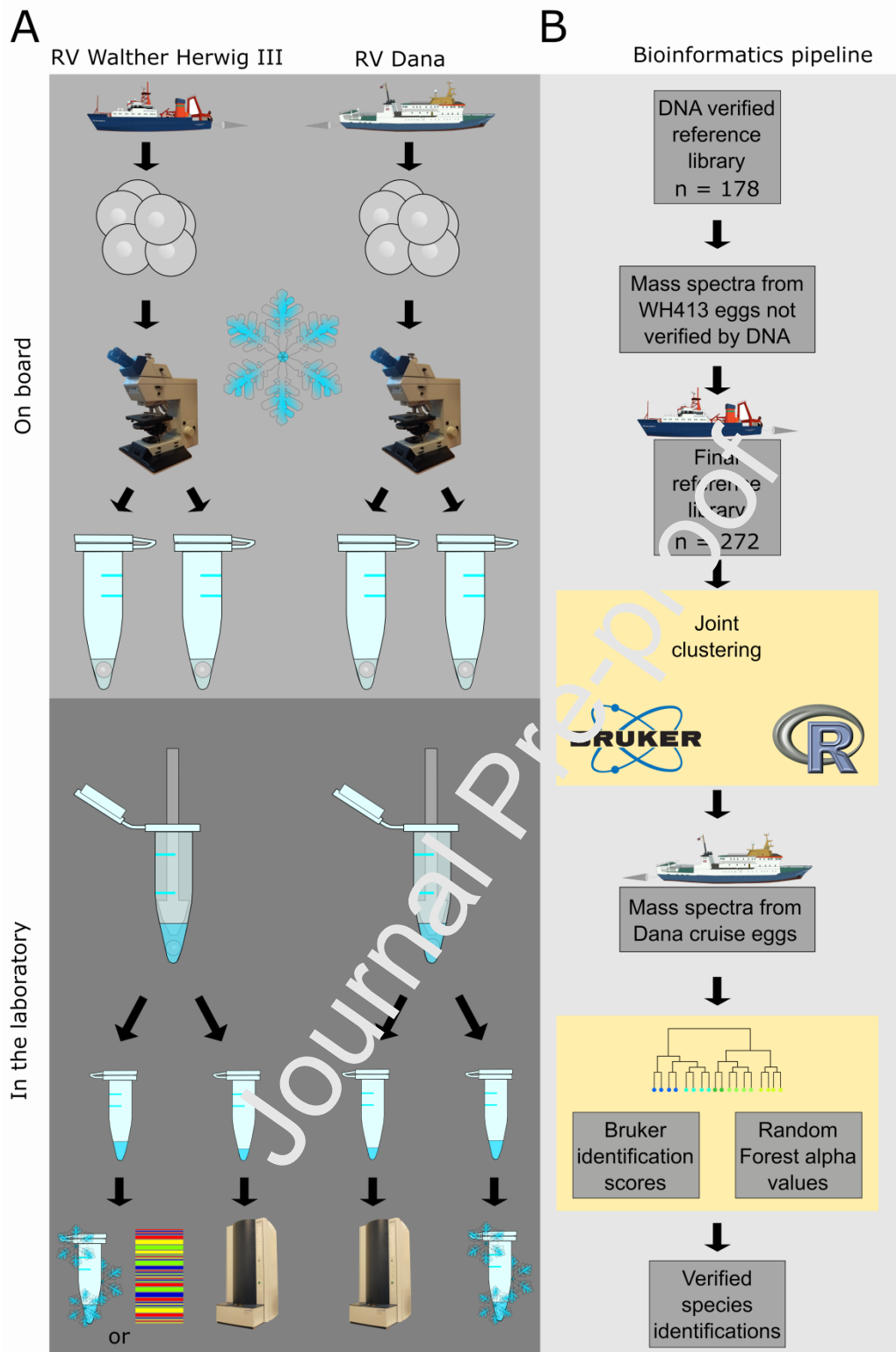


Fig. 2: **A** Sample processing on board was tried to be carried out as cooled as possible. After routine egg processing, eggs were separated into individual microcentrifuge tubes. In the laboratory these samples were further prepared by crushing with a pestle in 40 μ l molecular grade water. Of this mixture, 35 μ l were used for DNA extraction and 5 μ l for MALDI-TOF MS. If no DNA extraction was carried out, samples were stored at -20°C . **B** A DNA verified reference library based on samples from Walther Herwig III cruise samples was used to identify remaining WH413 eggs to construct a

final reference library for identification of samples from RV Dana by a line of evidence from different methods.

MALDI-TOF MS

Of the 40 μl mixture containing water and the crushed eggs, for all samples 5 μl were used for extraction of peptides and proteins (Fig. 2). To this mixture, 5 μl of 5% trifluoroacetic acid were added and incubated for 1h at room temperature (RT). After 1h, 1 μl of this solution was transferred to a MALDI-TOF target plate and dried at RT. The dried extract was covered with a layer of 1 μl α -Cyano-4-hydroxycinnamic acid (HCCA) as a saturated solution in 50% acetonitrile, 47.5% molecular grade water and 2.5% trifluoroacetic acid. For each egg, the solution was applied to one spot and air dried for co-crystallization of proteins and peptides with the matrix.

Measurements were carried out on a Microflex LT/SH System (Bruker Daltonics) using method MBTAuto with laser intensity between 50% and 60%. Peak evaluation during measurement was carried out in a mass peak range between 2k – 10k Dalton (Da) using a centroid peak detection algorithm, a signal to noise threshold of 2 and a minimum intensity threshold of 600, with a peak resolution higher than 400. Proteins/Oligonucleotide method was employed for fuzzy control with a maximal resolution ten times above the threshold. To create a sum spectrum, 240 satisfactory shots were summed up. Each spot was measured once. Before measurements, system calibration was carried out according to the Bruker default procedure using Bacterial standard (BTS).

Measurements for the reference library were carried out at different laboratories, but on same instrument brand and model. Two different instruments were used and compared. Reference specimens from plates TIP1 and TIP2 (WH413) were measured at the laboratory of the Lower Saxony State Office for Consumer Protection and Food safety in Cuxhaven (Cux) (Germany). Reference specimens from plates TIP3 and TIP4 (WH413) and TIP5, TIP6 and TIP12 (Dana cruise) were measured at the proteome laboratory of Senckenberg am Meer, German Centre for Marine Biodiversity Research (DZMB) in Wilhelmshaven (Whv).

Reference library preparation

For the reference library, 359 samples from cruise WH413 were distributed on four measurement plates: TIP1 (n=94), TIP2 (n=95), TIP3 (n=95) and TIP4 (n=75).

After measurements, an initial clustering of data was carried out and some specimens from different clusters were selected for DNA amplification. Thus, for 210 samples a COI barcode was assessed next to a MALDI-TOF mass spectrum. The mass spectra of these 210 samples were controlled separately

by eye for quality (e.g. low intensities or obvious signal degeneration) and low quality samples discarded. Further quality control was carried out in R. Mass spectra that were found with an exceptionally divergent A-score using quality control from the R-package MALDIrppa [52], as well as specimens with less than 35 mass peaks at SNR = 5 were discarded. This resulted in a mass spectra library for 178 samples from 10 species verified by DNA barcoding. Among the analyzed eggs the following species were represented: *Gadus morhua* Linnaeus, 1758, *Glyptocephalus cynoglossus* (Linnaeus, 1758), *Hippoglossoides platessoides* (Fabricius, 1780), *Lepidorhombus whiffiagonis* (Walbaum, 1792), *Limanda limanda* (Linnaeus, 1758), *Melanogrammus aeglefinus* (Linnaeus, 1758), *Merlangius merlangus* (Linnaeus, 1758), *Pleuronectes platessa* Linnaeus, 1758, *Pollachius virens* (Linnaeus, 1758) and *Trisopterus esmarkii* (Nilsson, 1855) (Supplementary Fig. 1).

During library preparation, a systematic mass shift between instruments in different laboratories was detected. Because this shift made aligning homologous mass peaks impossible, masses of measurements from Cuxhaven were adjusted by adding 7 Da to all measured masses.

To expand the library beyond samples previously identified by DNA, thus covering a higher mass spectra variability per species, the reference library was used to identify the remaining specimens from plates TIP1 to TIP4 by hierarchical clustering and Random Forest (RF) [53] with the *post-hoc* test described by Rossel & Martínez Arbizu (2015) [54] from the R package RFtools (<https://github.com/pmartinezarbizu/RFtools>) [55] using a 1% alpha value for false positive recognition. False positives were discarded, resulting in a data set containing 272 samples.

Identification of samples from RV Dana in R

In total, 243 eggs from RV Dana were measured, resulting in 239 successful measurements. Quality of mass spectra was checked using quality control from R-package MALDIrppa [52] by the command 'screenSpectra' (thScale=2.5; it=105; SigNoi=7; hws=20; tol=0.001). However, we deviated from the recommendation of applying the quality control to raw signal as we found applying it to already processed data (transformed intensity, smoothed, baseline corrected and normalized signal) identified low quality mass spectra more reliable. The A-score threshold for evaluation of mass spectra quality was also adjusted based on comparison to reference library quality A-scores (Fig. 3A). To assess quality, data was trimmed to a range between 5,000 to 20,000 Da. This was done because all mass spectra showed strong intensity signals between 2,000 to 5,000 Da even when mass spectra quality was in fact already degenerating (Fig. 3 C, D). Degenerated signals showed hardly any signal in the size range from 5,000 Da and higher. Trying to identify these mass spectra would result in unreliable species identification. Discarding of samples with a high A-score was further supported by low numbers of peaks (Fig. 3B). The majority of low quality mass spectra were found with less than 35 peaks, resulting in a bimodal histogram when regarding peak numbers of all samples (Fig. 3B). Characteristics such as signal intensities were not informative on mass spectra quality (Fig. 3C, D)

because low quality mass spectra still showed high signal intensities in a m/z range between 2,000 to 5,000 Da.

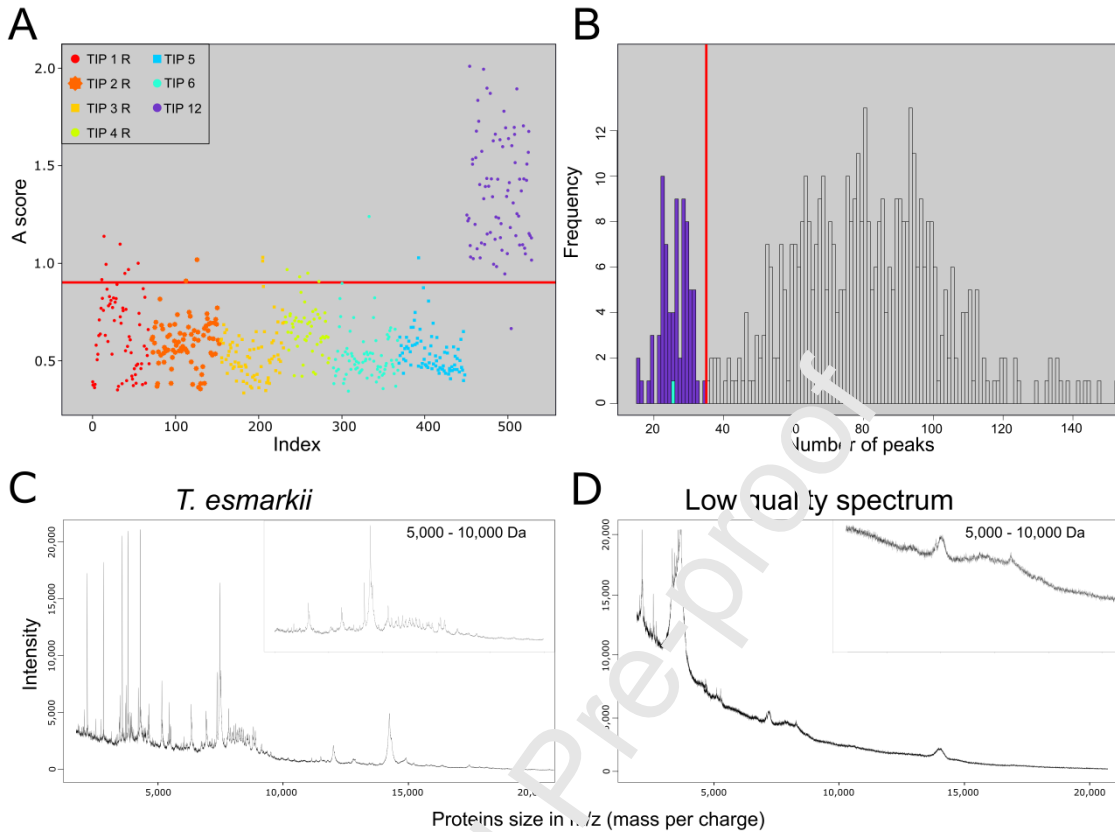


Fig. 3: **A** A-score plot from R-package `MAQDIrppa`. The A-score was calculated based on mass spectra trimmed to a size range from 5k to 20k Da. This was done because quality control failed to recognize all low quality mass spectra, when applied to the full mass range. **B** The histogram displays the number of peaks for all measured samples including the reference library. The red line is drawn at 35 peaks. Almost all mass spectra with less than 35 detected peaks were obtained from TIP12 (purple). **C** Mass spectrum of a good quality measurement and **D** a low quality mass spectrum. Even though high intensity signal was measured in a range between 2k to 5k Da, the range above 5k barely shows any good signal. In a good quality mass spectrum however, also in this size range, there are still some distinct peaks.

After excluding mass spectra of inferior quality, data was trimmed to an identical range from 2,000 to 20,000 Da. Intensities were square root transformed, and smoothed using Savitzky-Golay method [56]. Baseline correction was carried out with SNIP method [57] and intensities were calibrated with TIC method. Peak detection was done with a SNR of 5. Peak binning was carried out repeatedly until data set peak number decrease saturated. Finally, data was Hellinger transformed [58] for further analysis in R [54].

In order to obtain reliable identification in R, classification was based on different approaches. On the one hand, hierarchical clustering using Euclidean distances (Supplementary Fig. 2) and Ward's D cluster algorithm including reference spectra was applied. On the other hand identification was based on a RF classification approach (ntree=2,000; mtry=35; sampsize=6). To increase the power of the RF model, the sampsize was limited. This resulted in discarding of some species for the RF analyses. Overlapping identifications from both approaches were regarded as correct identifications. For the RF *post-hoc* test, we followed the recommendation of Rossel & Martínez Arbizu (2018) [59] to use a 1% α value additionally to a 5% α value.

Identification of samples from RV Dana using Bruker Biotyper®

The mass spectra used in R were also analyzed using the Bruker Biotyper®. From the reference library MSPs for each species were created using default settings (max. mass error of each single spectrum = 2,000; desired mass error for the MSP = 200; desired peak frequency minimum = 25%, max. desired peak number for the MSP = 70). Default settings for identification were used: frequency threshold for spectra adjusting = 50; frequency threshold for score calculation = 5; max. mass error of the raw spectrum = 2,000; des. Mass tolerance of the adjusted spectrum = 250; accepted mass tolerance of a peak = 600, parameter of the intensity correction function = 0.25. Mass spectra for identification were pre-processed by trimming mass spectra range from 2,000 to 20,000, smoothing using the Savitzky-Golay method with a frame size of 25 Da, baseline subtraction using the multipolygon method with a search window of 5 Da in two runs. Method applied for normalization was maximum norm and peak picking was done using spectra differentiation method with max. 100 peaks, a threshold of 0.001 and a SNR of 5. Default Biotyper® identification score thresholds to classify identification results were used (0 – 1.699 no reliable identification, 1.700-1.999 probable (genus) identification and ≥ 2.000 reliable identification) to classify identifications.

Identification agreement

To test the identification agreement of the different applied identification methods, percentage agreement was calculated for all methods conjoint using the command `agree` from the R package `irr` [60]. Classifications from the method to be tested against the remaining methods were used after application of a post-test. Thus, identifications failing the post-test were treated as unknowns. These were compared to raw classifications of the other methods.

Results

In order to construct a reference library for identification of fish eggs using MALDI-TOF MS confirmed by identification through DNA barcoding, specimens were chosen from an initial hierarchical clustering analysis. Of 359 specimens measured by two different Microflex LT/SH Systems, DNA barcoding was applied to these 210 samples whereas only 178 were used in the final reference library (Fig. 4).

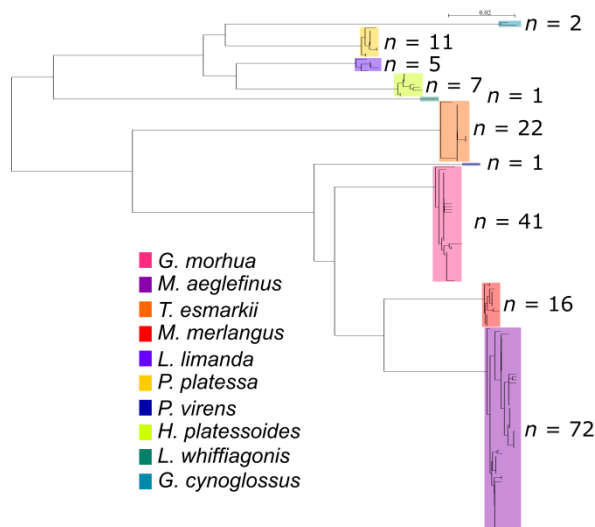


Fig. 4: NJ tree based on K2P distances of 178 mt COI sequences amplified from fish eggs. The tree displays the COI sequences for the samples used as reference specimens for the mass spectra library.

Between mass spectra measured on different instruments (TIP1+2 = Cux; TIP3+4 = Whv) a mass shift of 5 to 10 Da was detected (Fig. 5 A). Because peak binning or spectra alignment using reference peaks did not solve this problem, peak masses were manually adjusted by adding 7 Da to all masses measured in Cuxhaven. This resulted in improved comparability and binning of mass peaks from the two different instruments (Fig. 5 B). In general, the instrument used in Cuxhaven showed higher sensitivity for detection of molecules compared to the instrument employed in Wilhelmshaven. For instance, measurements of *M. aeglefinus* were affected by this difference in sensitivity by an average difference of 10 detected peaks between instruments (Cux = 96.4, n = 60; Whv = 86.0, n = 44). The same was found for *G. morhua* with 72.2 peaks on average measured in Wilhelmshaven (n = 43) whereas in Cuxhaven on average 97.6 Peaks (n = 18) were measured using the same instrument settings.

Not correcting the mass shift resulted in measurement-site specific clusters for several species within an initial clustering analysis (Fig. 5 C, indicated by stars). When checking the reference library for application in a RF classification approach, the RF model showed misclassification rate of 0.56%, probably caused by this mass shift. The misidentification rate was reduced to zero after adjusting the masses. Furthermore, the clusters then also referred to COI based species identification (Fig. 5 D)

rather than to instruments used. This DNA verified and quality curated reference library was then used to identify the remaining samples from cruise WH413. The final library additionally containing specimens assigned to species using the reference library, consisting of samples from this cruise, contained a total of 272 specimens from ten species (*G. cynoglossus*: n = 4; *G. morhua*: n = 60; *H. platessoides*: n = 11; *L. limanda*: n = 6; *L. whiffiagonis*: n = 3; *M. aeglefinus*: n = 104; *M. merlangus*: n = 25; *P. platessa*: n = 29; *P. virens*: n = 1; *T. esmarkii*: n = 30). This extended data set showed a misidentification rate of 0 as well and was used for further identification of samples from RV Dana cruise 02/2018 using Bruker Biotyper® software, hierarchical clustering in R and the RF approach with a *post-hoc* test for false positive recognition.

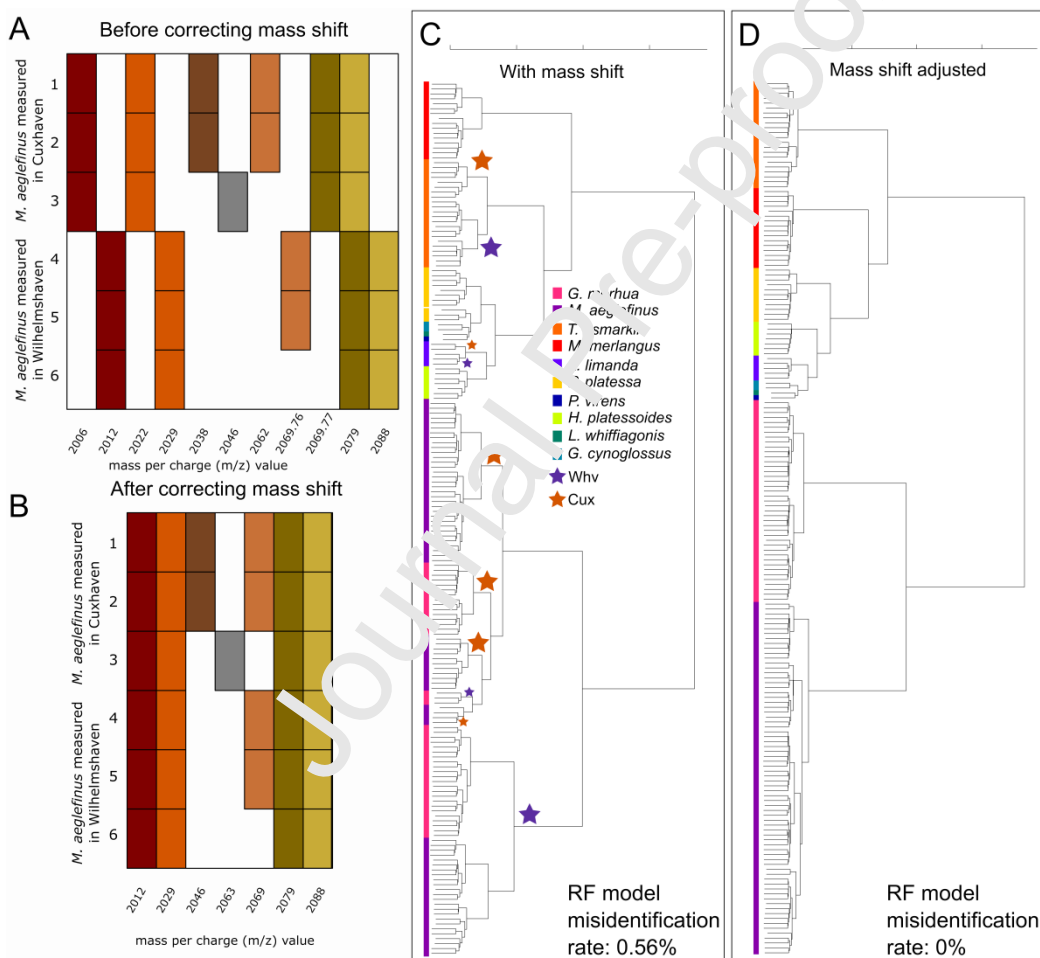


Fig. 5: A mass shift occurred between instruments from Cuxhaven (specimens 1-3) and Wilhelmshaven (specimens 4-6) displayed by comparing mass spectra of six *M. aeglefinus* samples in the m/z range between 2,000 and 2,090 Da. Each rectangle displays a peak. The m/z value is displayed after peak binning. In **A** the peaks are displayed before manually correcting the mass shift. In **B** the peaks are displayed after manually correcting the mass shift. **C** Hierarchical cluster analyses of the quality controlled, extended data set containing DNA verified specimens and specimens from cruise

WH413 identified using the reference library prior to correction of the mass shift between instruments and **D** after correcting it. For several species, this mass shift impeded species-specific clustering and fostered instrument-specific clusters (marked by stars).

Identification of the remaining 158 mass spectra from Dana cruise after quality control was carried out using a line of evidence from three different approaches (Fig. 2). I) Identification through hierarchical clustering with reference specimens. II) Identification using the Bruker Biotyper® software based on MSPs generated from DNA verified mass spectra. III) Identification using RF based on a reference library generated from DNA verified mass spectra (Fig. 6).

Concordant identifications were found for three samples as *G. morhua* and 151 as *P. platessa*. One specimen was identified by clustering and by the Bruker Biotyper® as *M. virens*. This species was due to the minimum number of specimens per species in the model ($n=1$) not included in the RF library and the specimen was thus identified by RF as *M. merlangus*. Three specimens identified by clustering and RF as *P. platessa* were identified by the Bruker Biotyper® as *G. cynoglossus* (Fig. 6, black bars). Identifications that were not concordant between the different methods were recognized either as unreliable identification by the Biotyper® or as false positive by the *post-hoc* test. Even though all three methods resulted in concordant results for most of the species, some classifications were rejected by post-tests using default settings. Of the 157 concordant identifications, the Biotyper® recognized 52 as correct identifications, 43 as probable (genus) identifications and 63 as not reliable identifications. Only one of the concordant RF classifications was recognized as false positive by the 1% alpha of the RF *post-hoc* test (Fig. 6).

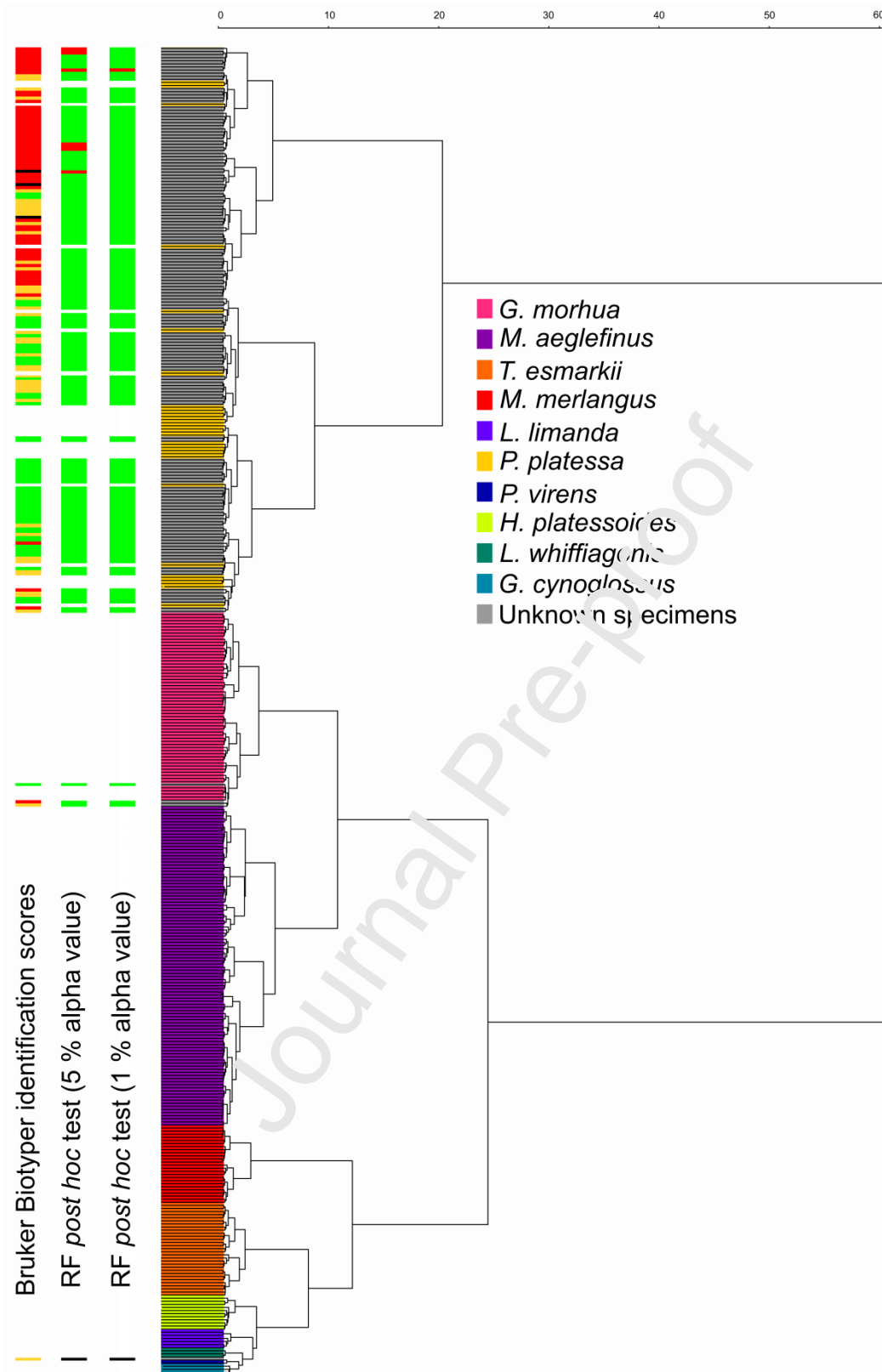


Fig. 6: Hierarchical clustering of unknown mass spectra and reference library. Colors at the tip of the tree indicate species affiliation of reference specimens. Grey bars indicate unknowns. Bars beside the tree indicate classification quality using default thresholds from the Biotyper® software (most left) and different α values for RF classification *post-hoc* test: 5% (middle), 1% (most right). Green marks

true positives, orange probable identifications and red false positives. Black bars mark identifications deviating between the different methods identifications.

Discussion

MALDI-TOF MS for rapid species identification

With this study we demonstrated that MALDI-TOF MS is a suitable tool for identification of fish eggs. We found that groups based on mass spectra are congruent with clusters based on DNA barcoding [42,44,61,62]. MALDI-TOF MS can thus be considered a versatile alternative to expensive COI barcoding [43] and difficult morphological identification of fish eggs [10] in ecological studies or fish surveys. However, a well curated and extensive reference library is necessary for reliable results and species detection. Even though a comprehensive reference library would be desirable, previous studies have shown that, for instance using RF classification, even with an incomplete library, detection of unknown species is possible [43,54]. Hence, enabling surveys or biodiversity assessments also in areas for which a complete library is not available yet. If potential new species are detected by MALDI-TOF MS, our sample preparation would allow an additional DNA analysis if DNA extract was stored properly. Moreover using a line of evidence from different classification methods in an integrative approach reinforces identifications and improves false positive detection.

The advantage of using a line of evidence for identification and the necessity of comprehensive reference libraries is especially emphasized by the rejection of several *P. platessa* identifications by the Bruker Biotyper®. Regarding the positions of the rejected identifications in the hierarchical clustering analysis (Fig. 6), it seems that the rejections often depend on an incomplete coverage of mass spectra variability in the reference library. Within the clustering analysis, the species *P. platessa* is divided into two main groups of which some clusters include no reference specimens at all. Whereas the majority of unknown specimens are situated within the upper group, only four reference samples are found therein. While the Bruker Biotyper® seems to have difficulties working with this lack of data, the RF approach shows similar identifications which are, in opposite to the Bruker identifications, widely supported by the *post-hoc* test. However, by combining the results of the different methods, identifications can be accepted more confidently. Mere clustering for instance is prone to misidentifications because of unrecognized species margins or missing reference specimens clustering with specimens to be identified. In particular those clusters consisting only of few specimens can hardly be resolved for reliable identifications.

We found quality control to play an important role for instance to recognize probable systematic mass-shifts between instruments or low quality mass spectra. Trying to identify irregular mass spectra will probably lead to misidentifications. Here, the screenSpectra function from the R-package MALDirppa

was, after some adjustments of settings, very helpful in automatically detecting low quality mass spectra. This was additionally supported by extraordinary low peak numbers recognizable due to a bimodal distribution of peak abundances. As was reported before by several authors, good sample preservation is of highest importance for successful application of this method [59,63–65].

This is also supported by our results for samples from a single cruise that differed highly in mass spectra quality. Samples measured on TIP5 and TIP6 belonged to eleven different stations sampled with RV Dana with a low number of eggs per site that were sorted at a time. Generally, samples were preserved only after all samples from a certain station were processed, which here also included taking images. When processing only few samples, eggs were quickly put into preservation. However, when numerous eggs were processed at a certain station, also the eggs were exposed to adverse conditions for longer times. Presumably, this is what affected eggs measured on TIP12. In total, 114 eggs were sampled on this station, prolonging the processing time, thus causing the poor mass spectra qualities. For future application, it needs clear sampling protocols which include limitation of handling times for each individual specimen.

A crucial step to a wider application of this method in fishery surveys or biological and ecological studies would be the creation of a publicly accessible database for MALDI-TOF mass spectra comparable to BOLD [66] or GenBank [67] for DNA applications. Even though many mass spectrometry data bases already exist, at the moment none of these really complies with the kind of data produced here. That is why data is often not made publicly available or deposited in unspecific data repositories such as Dryad Data repository.

Inter-laboratory application and identification thresholds

Additionally to the good identification success, we were able to show that data resulting from different instruments (same brand and model) can be made interoperable between laboratories, thus allowing for species identifications using reference libraries produced elsewhere and do not only need to rely on in-house databases. Even though our test of inter-laboratory compliance of measured proteome fingerprint only involved two instruments, it gives evidence for the general usability of MS data libraries from different studies of metazoans provided that specimens were prepared identically. For microorganisms such as bacteria or fungi this was shown several times before [68–70]. For species identification of metazoan however, we were unable to find a study supporting inter-laboratory use of MS data. Most studies mention the use of in-house reference libraries [71]. Nonetheless, inter-laboratory data has to be used with caution, as different instruments produced mass spectra with different detection sensitivities, probably resulting in instrument or laboratory specific signal [69]. In our study, a mass shift between instruments occurred but was noticed due to an initial clustering of specimens for which the species affiliation was verified using DNA barcoding. The shift could be visualized using the R-package MALDIrppa [52] (Fig. 5) and thus corrected manually. If future

development reaches towards a wide implementation of MALDI-TOF MS for species identification of metazoans, general databases will be necessary. Mass shifts such as the one that occurred in the present study could then easily be recognized by the use of a defined internal size standard to improve data alignment. Even though clustering helped detecting the mass shift, trying to identify the specimens measured with one instrument based on the reference from the other instrument before correction of the mass shift by clustering would have failed for the majority of the specimens as is emphasized by figure 5 C. Here, the species *G. morhua* and *M. aeglefinus* for instance would not be separable by mere clustering.

For the two supervised methods, we visualized the impact of the mass shift on the identification process in figure 7, comparing identification values and thresholds from different set ups for the species *P. platessa* including also reference specimens in the identification process. Whereas the RF classification based on reference specimens from a single instrument widely remains the same compared to a complete reference library, the *post-hoc* results change clearly (Fig. 7 A, B, D, E). Using the Cux reference library without mass shift adjustment results for the concordant *P. platessa* identifications in 69 classifications recognized as false positives by the *post-hoc* test (Fig. 7 B, orange line). In contrast to this, only 27 are regarded as false positives after adjusting the mass shift (Fig. 7A, orange line). Because the majority of mass spectra were measured in Whv, the RF classification based on Whv specimens seems not to be influenced as much. None of the concordant *P. platessa* identifications is recognized as false positive with the 1% threshold (Fig. 7 D, E, orange line). Nevertheless, RF values for the Cux animals are remarkably lower without mass shift correction (Fig. 7 E). This effect however diminishes when reference libraries are combined (Fig. 7 G, H). Here, not adjusting the mass shift even positively influenced the identification success (Fig. 7 H). One of the reference specimens that showed low RF values in all analyses now is accepted as a correct identification under the 1% alpha value.

In general, the applied alpha values for false positive recognition seem to be well chosen. In approaches including the Whv specimens in the reference to create a RF model, the majority of concordant identifications are accepted as correct identifications under the 1% alpha value. However, when Cux specimens were not included in the model, up to two specimens from species not included in the RF model, were accepted as correct identifications (Fig. 7 D, E).

For the Bruker Biotyper® identification it was not possible to simply adjust the mass shift. Specimens were thus identified without correction, however testing the different reference libraries as well (Fig. 7 C, F, I). Generally, the Biotyper® thresholds are stricter than the chosen alpha values for the *post-hoc* test. With the Cux library, for 77 of the concordant *P. platessa* identifications an identification value of less than 1.7, which is the most relaxed Biotyper® default value, was calculated (Fig. 7 C, orange line). When the Whv library is used, still 40 classifications are regarded as incorrect (Fig. 7 F, orange line). Finally, combining the two reference libraries resulted in 59 specimens being recognized as

misclassifications (Fig. 7 I, orange line). So in contrast to RF, creating MSP from divergent reference mass spectra rather weakened the identification capabilities than increasing it.

Our results thus emphasize that even the most relaxed Biotyper® default threshold is still far stricter than α values applied for RF classifications. The thresholds may be well applicable and established for cultured microorganisms that show low biological variabilities but may be inadequate for identification of field caught metazoan specimens that exhibit higher mass spectra variability due to different environmental conditions and mass spectra variation due to inferior sample storage [59,63–65].

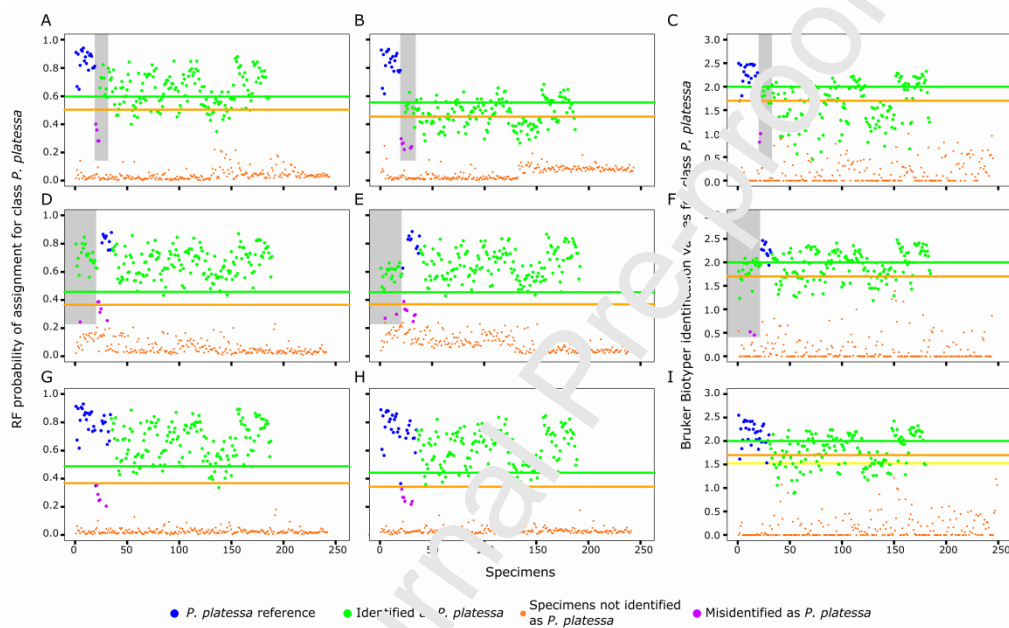


Fig. 7: Identification values for *P. platessa* from the different applied supervised identification methods. The right column contains identifications using the Bruker Biotyper® software. The other two columns depict RF results. **ADG**: Mass shift adjusted. **BEH** Mass shift not adjusted. In the upper row only specimens from Cux were used for reference. The second row shows identifications when only Whv specimens were used as library. In the lowest row, all library specimens were used for classification of unknowns but also the reference itself. Grey bars indicate the specimens from the reference library that were not used in the corresponding approach. Orange line: 1.0% α value for RF *post-hoc* test; 1.7 threshold for Biotyper® scores. Green line: 5.0% α value for RF *post-hoc* test; 2.0 threshold for Biotyper® scores. Yellow line: lowest Biotyper® value of a reference specimen.

Different studies on both microorganisms and metazoans have already shown that adjusting identification values based on known specimens might be necessary [71,72] to obtain good

identification success. Our data emphasizes the necessity of taking identification values of specimens from other species into account (Fig. 7 C, F, I: orange dots). Even though these specimens were identified as different species, they receive identification scores for all other species for which a MSP was included in the identification process. Sometimes, these scores can rise as high as values of specimens actually belonging to the identified species, thus creating the risk of producing false positives or ambiguous identifications. Concluding from this, lowering identification score values is necessary but should be carried out with care. Nonetheless, figure 7 I for instance shows lowering the threshold to the lowest value for library specimens would already result in 25.4% less unreliable identifications without causing false positives (Fig. 7 I, Fig. 8, yellow line). Based on our results we would however not recommend a general reduction of score values as this would not consider species specific differences.

When testing the agreement of the different classification methods while varying thresholds, we can see that lowering thresholds for both methods increases the identification agreement (Fig. 8) between all three methods. Nevertheless, it needs to be taken into account that by lowering these thresholds, also more misidentifications may be recognized as correct identifications. In our study, lowering the alpha value for the RF *post-hoc* test would increase the identification agreement marginally (Fig. 8) but at the same time almost completely eliminate the false positive recognition (Fig. 7 H).

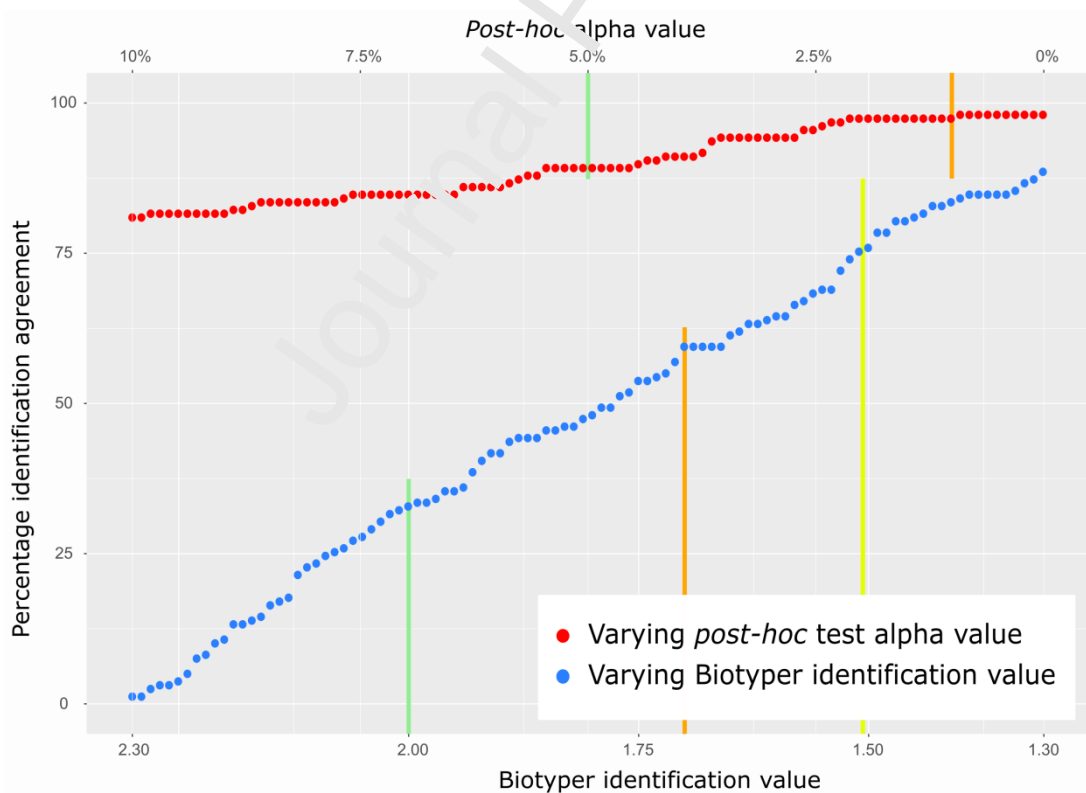


Fig. 8: Alteration of identification agreement between the three classification methods under changing thresholds of post-hoc test values (red) and Biolyper® values (blue). Lowering the identification values increases the agreement between the methods. Green lines display the strict thresholds and

orange lines the slacker thresholds. The yellow line marks the lowest Biotyper® value of a reference specimen.

Conclusion

Morphological identification of fish eggs is difficult and demands a lot taxonomic knowledge. However, in some cases differentiation even between lesser related species is not possible because of lacking morphological diagnostic features, particularly in eggs at early developmental stages. Thus, a rapid and inexpensive identification method such as MALDI-TOF MS can help to improve fish egg identification and maybe even accelerate work in large surveys where many specimens have to be identified. MALDI-TOF MS is a promising candidate for future fish egg monitoring. The method provides a reliable alternative to morphological or molecular identification. At the same time it is comparably cheap and demands only little sample preparation effort. Compared to similar approaches such as DNA-barcoding, results are available faster and can be analyzed quicker. Future applications need to aim at generating more comprehensive databases. Using a line of evidence for species identification increases identification confidence. By lowering identification thresholds, identification success can be increased. However, this needs to take into account, that lowering thresholds may result in an increase of false positives.

Data availability

Sequence data is stored in BOLD in the Project “FEM:Fish eggs identification using MALDI-TOF MS”. Data can be found on GenBank using accession numbers XXX-XXX. MALDI-TOF MS raw data and subsequent metadata is stored in Dryad data repository: DOIXXX.

ORCID author statement

Sven Rossel: Writing - Original Draft, Visualization, Data Curation, Formal analysis, Conceptualization **Andrea Marco:** Writing - Review & Editing, Investigation, Conceptualization **Matthias Kloppmann:** Resources, Funding acquisition, Writing - Review & Editing, Conceptualization **Pedro Martínez Arbizu:** Resources, Funding acquisition, Writing - Review & Editing **Bastian Huwer:** Resources, Writing - Review & Editing **Thomas Knebelsberger:** Writing - Review & Editing, Investigation, Conceptualization

The following are the supplementary data related to this article.

Supplementary Fig. 1: Mass spectra for all species included in our analysis. Even in a m/z range higher than 5,000 all species still show distinct signals. This allows using this part of the spectra for quality control because frequently, low quality mass spectra show no more signal in this area.

Supplementary Fig. 2: Matrix containing Euclidean distances of all specimens in the data set sorted by species. Dark colors indicate small distances. Light colors indicate larger distances. Whereas the distances between most of the species is high, the distances between the two gadid species *M. aeglefinus* and *G. morhua* are comparably small.

Acknowledgement

We would like to thank Maik Tiedemann and Sakis Kroupis, Thünen Institute of Sea Fisheries, for sorting MIKey-M net samples and meticulously measuring, staging and preserving the fish eggs during WH 413 cruise. Maik Tiedemann received funding for sampling on the International Bottom Trawl Survey through the European Maritime and Fisheries Fund of the European Union (Data Collection Framework). Furthermore, we highly appreciate the support and assistance with sampling by the officers and crews of the FRV Walther Herwig III and RV Dana during the International Bottom Trawl Survey (IBTS) Quarter 1 in 2018. We thank Dr. Ralf Pund and Marko Kranz from LAVES (Lower Saxony State Office for Consumer Protection and Food safety) Cuxhaven for their support and the possibility to use their MALDI-TOF instrument. This is publication no 11 of Senckenberg am Meer Proteome Laboratory.

Conflict of interest

Conflict of interest: biome-id offers molecular-based commercial services including DNA-barcoding and MALDI-TOF analysis to monitoring agencies and research institutes.

CRedit author statement

Sven Rossel: Writing - Original Draft, Visualization, Data Curation, Formal analysis, Conceptualization **Andrea Barco:** Writing - Review & Editing, Investigation, Conceptualization **Matthias Kloppmann:** Resources, Funding acquisition, Writing - Review & Editing, Conceptualization **Pedro Martínez Arbizu:** Resources, Funding acquisition, Writing - Review & Editing **Bastian Huwer:** Resources, Writing - Review & Editing **Thomas Knebelsberger:** Writing - Review & Editing, Investigation, Conceptualization

References

- [1] Fox, C.J., Taylor, M., Dickey-Collas, M., Fossum, P., et al., Mapping the spawning grounds of North Sea cod (*Gadus morhua*) by direct and indirect means. *P ROY SOC B-BIOL SCI* 2008, 275, 1543–1548.
- [2] Ibaibarriaga, L., Irigoien, X., Santos, M., Motos, L., et al., Egg and larval distributions of seven fish species in north-east Atlantic waters. *Fish. Oceanogr.* 2007, 16, 284–293.
- [3] Ahern, A.L.M., Gómez-Gutiérrez, J., Aburto-Oropeza, O., Saldierna-Martínez, R.J., et al., DNA sequencing of fish eggs and larvae reveals high species diversity and seasonal changes in spawning activity in the southeastern Gulf of California. *Mar. Ecol. Prog. Ser.* 2018, 592, 159–179.
- [4] Harada, A.E., Lindgren, E.A., Hermsmeier, M.C., Rogowski, P.A., et al., Monitoring spawning activity in a southern California marine protected area using molecular identification of fish eggs. *PloS one* 2015, 10, e0134647.
- [5] Lockwood, S., Nichols, J., Dawson, W.A., The estimation of a mackerel (*Scomber scombrus* L.) spawning stock size by plankton survey. *J. Plankton Res.* 1981, 3, 217–233.
- [6] Moser, H.G., Charter, R.L., Watson, W., Ambrose, D., et al., The CalCOFI ichthyoplankton time series: potential contributions to the management of rocky-shore fishes. *CAL COOP OCEAN FISH* 2001, 112–128.
- [7] Armstrong, M., Connolly, P., Nash, R., Pawson, M., et al., An application of the annual egg production method to estimate the spawning bio mass of cod (*Gadus morhua* L.), plaice (*Pleuronectes platessa* L.) and sole (*Solea solea* L.) in the Irish Sea. *ICES J MAR SCI* 2001, 58, 183–203.
- [8] Köster, F.W., Huwer, B., Kraus, G., Diekmann, R., et al., Egg production methods applied to Eastern Baltic cod provide indices of spawning stock dynamics. *Fish. Res* 2020, 227, 105553.
- [9] Koslow, J.A., Wright, M., Ichthyoplankton sampling design to monitor marine fish populations and communities. *Mar. Policy* 2016, 68, 55–64.
- [10] Russell, F.S., *The eggs and planktonic stages of British marine fishes*, vol. 524, Academic press London, 1976.
- [11] Munk, P., Nielsen, J.G., *Eggs and larvae of North Sea fishes*, Biofolia, Frederiksberg, Denmark 2005.
- [12] Hempel, G., *Early life history of marine fish: the egg stage*, Washington Sea Grant WA, USA, 1979.
- [13] ICES, Manual for the egg survey for winter spawning fish in the North Sea. *Series of ICES Survey Protocols SISP 13* 2018, 19.
- [14] ICES, Manual for mackerel and horse mackerel egg surveys, sampling at sea. *Series of ICES Survey Protocols SISP 6* 2019, 82.
- [15] ICES, Report of the Workshop on egg staging, fecundity, and atresia in horse mackerel and mackerel (WKFATHOM2). 8-12 October and 19-23 November. Bremerhaven, Germany and IJmuiden, Netherlands. *ICES CM 2018/EOSG:22* 2018, 74.
- [16] Pappalardo, A.M., Petracchioli, A., Capriglione, T., Ferrito, V., From fish eggs to fish name: Caviar species discrimination by COI Bar-RFLP, an efficient molecular approach to detect fraud in the caviar trade. *Molecules* 2019, 24, 2468.
- [17] Choi, H., Oh, J., Kim, S., Genetic identification of eggs from four species of Ophichthidae and Congridae (Anguilliformes) in the northern East China Sea. *PloS one* 2018, 13, e0195382.
- [18] Hofmann, T., Knebelberger, T., Kloppmann, M., Ulleweit, J., Raupach, M., Egg identification of three economical important fish species using DNA barcoding in comparison to a morphological determination. *J APPL ICHTHYOL* 2017, 33, 925–932.

- [19] Taylor, M.I., Fox, C., Rico, I., Rico, C., Species-specific TaqMan probes for simultaneous identification of (*Gadus morhua* L.), haddock (*Melanogrammus aeglefinus* L.) and whiting (*Merlangius merlangus* L.). *Molecular Ecology Notes* 2002, 2, 599–601.
- [20] Welker, M., Proteomics for routine identification of microorganisms. *Proteomics* 2011, 11, 3143–3153.
- [21] Singhal, N., Kumar, M., Kanaujia, P.K., Viridi, J.S., MALDI-TOF mass spectrometry: an emerging technology for microbial identification and diagnosis. *FRONT MICROBIOL* 2015, 6.
- [22] Barbuddhe, S.B., Maier, T., Schwarz, G., Kostrzewa, M., et al., Rapid identification and typing of listeria species by matrix-assisted laser desorption ionization-time of flight mass spectrometry. *Applied and environmental microbiology* 2008, 74, 5402–5407.
- [23] Nagy, E., Maier, T., Urban, E., Terhes, G., et al., Species identification of clinical isolates of *Bacteroides* by matrix-assisted laser-desorption/ionization time-of-flight mass spectrometry. *Clinical Microbiology and Infection* 2009, 15, 796–802.
- [24] Calderaro, A., Arcangeletti, M.-C., Rodighiero, I., Buttrini, M., et al., Matrix-assisted laser desorption/ionization time-of-flight (MALDI-TOF) mass spectrometry applied to virus identification. *Sci. Rep.* 2014, 4, 6803.
- [25] La Scola, B., Campocasso, A., N'Dong, R., Fournier, G., et al., Tentative characterization of new environmental giant viruses by MALDI-TOF mass spectrometry. *Intervirology* 2010, 53, 344–353.
- [26] Chalupová, J., Raus, M., Sedlářová, M., Šebela, M., Identification of fungal microorganisms by MALDI-TOF mass spectrometry. *Biotechnol. Adv.* 2014, 32, 230–241.
- [27] Danezis, G.P., Tsagkaris, A.S., Camina, F., Brusic, V., Georgiou, C.A., Food authentication: Techniques, trends & emerging approaches. *TRAC-TREND ANAL CHEM* 2016, 85, 123–132.
- [28] Volta, P., Riccardi, N., Lauceri, K., Tonolla, M., Discrimination of freshwater fish species by Matrix-Assisted Laser Desorption/Ionization-Time Of Flight Mass Spectrometry (MALDI-TOF MS): a pilot study. *J LIMNOL* 2012, 71, e17.
- [29] Mazzeo, M.F., Giulio, B.D., Guerriero, G., Ciarcia, G., et al., Fish authentication by MALDI-TOF mass spectrometry. *J. Agric. Food Chem.* 2008, 56, 11071–11076.
- [30] Mazzeo, M.F., Siciliano, P.A., Proteomics for the authentication of fish species. *J PROTEOMICS* 2016, 17, 119–124.
- [31] Spielmann, G., Huber, I., Maggipinto, M., Haszprunar, G., et al., Comparison of five preparatory protocols for fish species identification using MALDI-TOF MS. *EUR FOOD RES TECHNOL* 2018, 244, 685–694.
- [32] Stahl, A., Schröder, U., Development of a MALDI-TOF MS-based protein fingerprint database of common food fish allowing fast and reliable identification of fraud and substitution. *Journal of agricultural and food chemistry* 2017, 65, 7519–7527.
- [33] Bi, H., Zhong, C., Shao, M., Wang, C., et al., Differentiation and Authentication of Fishes at Species Level Through Analysis of Fish Skin by MALDI TOF MS. *RAPID COMMUN MASS SP* 2019.
- [34] Salla, V., Murray, K.K., Matrix-assisted laser desorption ionization mass spectrometry for identification of shrimp. *Analytica chimica acta* 2013, 794, 55–59.
- [35] Maász, G., Takács, P., Boda, P., Várbiró, G., Pirger, Z., Mayfly and fish species identification and sex determination in bleak (*Alburnus alburnus*) by MALDI-TOF mass spectrometry. *Science of The Total Environment* 2017, 601, 317–325.
- [36] Flaudrops, C., Armstrong, N., Raoult, D., Chabrière, E., Determination of the animal origin of meat and gelatin by MALDI-TOF-MS. *J Food Compost Anal* 2015, 41, 104–112.

- [37] Sassi, M., Arena, S., Scaloni, A., MALDI-TOF-MS platform for integrated proteomic and peptidomic profiling of milk samples allows rapid detection of food adulterations. *J. Agric. Food Chem.* 2015, 63, 6157–6171.
- [38] Raharimalala, F., Andrianinarivomanana, T., Rakotondrasoa, A., Collard, J., Boyer, S., Usefulness and accuracy of MALDI-TOF mass spectrometry as a supplementary tool to identify mosquito vector species and to invest in development of international database. *MED VET ENTOMOL* 2017, 31, 289–298.
- [39] Loaiza, J.R., Almanza, A., Rojas, J.C., Mejia, L., et al., Application of matrix-assisted laser desorption/ionization mass spectrometry to identify species of Neotropical Anopheles vectors of malaria. *Malar. J.* 2019, 18, 95.
- [40] Yssouf, A., Flaudrops, C., Drali, R., Kernif, T., et al., Matrix-assisted laser desorption ionization-time of flight mass spectrometry for rapid identification of tick vectors. *J CLIN MICROBIOL* 2013, 51, 522–528.
- [41] Mathis, A., Depaquit, J., Dvovrák, V., Tuten, H., et al., Identification of phlebotomine sand flies using one MALDI-TOF MS reference database and two mass spectrometer systems. *PARASITE VECTOR* 2015, 8, 266.
- [42] Bode, M., Laakmann, S., Kaiser, P., Hagen, W., et al., Travelling diversity of deep-sea copepods using integrated morphological and molecular techniques. *J. Plankton Res.* 2017, 39, 600–617.
- [43] Rossel, S., Khodami, S., Martínez Arbizu, P., Comparison of rapid biodiversity assessment of meiobenthos using MALDI-TOF MS and Metabarcoding. *Front. Mar. Sci.* 2019, 6, 659.
- [44] Kaiser, P., Bode, M., Cornils, A., Hagen, W., et al., High-resolution community analysis of deep-sea copepods using MALDI-TOF protein fingerprinting. *DEEP-SEA RES PT I* 2018.
- [45] ICES, Manual for the Midwater Ring Net sampling during IBTS Q1. *Series of ICES Survey Protocols SISP 2* 2017, 25.
- [46] Ivanova, N.V., Zemlak, T.S., Parker, R.H., Hebert, P.D., Universal primer cocktails for fish DNA barcoding. *Mol. Ecol. Notes* 2007, 7, 544–548.
- [47] Ward, R.D., Zemlak, T.S., Jones, B.H., Last, P.R., Hebert, P.D.N., DNA barcoding Australia's fish species. *Philos. Trans. R. Soc. Lond., B, Biol. Sci.* 2005, 360, 1847–57.
- [48] Gouy, M., Guindon, S., Gascuel, O., SeaView version 4: a multiplatform graphical user interface for sequence alignment and phylogenetic tree building. *Mol. Biol. Evol.* 2010, 27, 221–224.
- [49] Edgar, R.C., MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.* 2004, 32, 1792–1797.
- [50] Saitou, N., Nei, M., The neighbor-joining method: a new method for reconstructing phylogenetic trees. *Mol. Biol. Evol.* 1987, 4, 406–25.
- [51] Kimura, M., A simple method for estimating evolutionary rates of base substitutions through comparative studies of nucleotide sequences. *J. Mol. Evol.* 1980, 16, 111–120.
- [52] Palarea-Albaladejo, J., Mclean, K., Wright, F., Smith, D.G., MALDIrppa: quality control and robust analysis for mass spectrometry data. *Bioinformatics* 2017, 34, 522–523.
- [53] Breimann, L., Random Forests. *Mach Learn.* 2001, 45, 5–32.
- [54] Rossel, S., Martínez Arbizu, P., Automatic specimen identification of Harpacticoids (Crustacea: Copepoda) using Random Forest and MALDI-TOF mass spectra, including a post hoc test for false positive discovery. *METHODS ECOL EVOL* 2018, 00, 1–14.
- [55] Martínez Arbizu, P., Rossel, S., RfTools: Miscellaneous Tools For Random Forest Models 2018.
- [56] Savitzky, A., Golay, M.J., Smoothing and differentiation of data by simplified least squares procedures. *Anal. Chem.* 1964, 36, 1627–1639.

- [57] Ryan, C., Clayton, E., Griffin, W., Sie, S., Cousens, D., SNIP, a statistics-sensitive background treatment for the quantitative analysis of PIXE spectra in geoscience applications. *Nucl. Instrum. Methods Phys. Res* 1988, 34, 396–402.
- [58] Legendre, P., Gallagher, E.D., Ecologically meaningful transformations for ordination of species data. *Oecologia* 2001, 129, 271–280.
- [59] Rossel, S., Martínez Arbizu, P., Effects of Sample Fixation on Specimen Identification in Biodiversity Assemblies based on Proteomic Data (MALDI-TOF). *Front. Mar. Sci.* 2018, 5, 149.
- [60] Gamer, M., Lemon, J., Singh, I.F.P., *irr: Various Coefficients of Interrater Reliability and Agreement*, 2019.
- [61] Rossel, S., Martínez Arbizu, P., Revealing higher than expected diversity of Harpacticoida (Crustacea: Copepoda) in the North Sea using MALDI-TOF MS and molecular barcoding. *Sci. Rep.* 2019, 9, 9182.
- [62] Riccardi, N., Lucini, L., Benagli, C., Welker, M., et al., Potential of matrix-assisted laser desorption/ionization time-of-flight mass spectrometry (MALDI-TOF MS) for the identification of freshwater zooplankton: a pilot study with three Eudiaptomus (Copepoda: Diaptomidae) species. *J. Plankton Res.* 2012, 34, 484–492.
- [63] Yssouf, A., Parola, P., Lindström, A., Lilja, T., et al., Identification of European mosquito species by MALDI-TOF MS. *Parasitol. Res.* 2014, 113, 2375–8.
- [64] Dieme, C., Yssouf, A., Vega-Rúa, A., Berenger J.-M., et al., Accurate identification of Culicidae at aquatic developmental stages by MALDI-TOF MS profiling. *PARASITE VECTOR* 2014, 7, 544.
- [65] Kaufmann, C., Ziegler, D., Schaffner, F., Crompter, S., et al., Evaluation of matrix-assisted laser desorption/ionization time-of-flight mass spectrometry for characterization of *Culicoides nubeculosus* biting midges. *MED VET ENTOMOL* 2011, 25, 32–38.
- [66] Ratnasingham, S., Hebert, P.D., BOLD: The Barcode of Life Data System (<http://www.barcodinglife.org>). *Mol. Ecol. Resour.* 2007, 7, 355–364.
- [67] Benson, D.A., Cavanaugh, M., Clark, K., Karsch-Mizrachi, I., et al., GenBank. *Nucleic Acids Res.* 2012, 41, D36–D42.
- [68] Vlek, A., Kolecka, A., Khashan, K., Theelen, B., et al., Interlaboratory comparison of sample preparation methods, database expansions, and cutoff values for identification of yeasts by matrix-assisted laser desorption ionization-time of flight mass spectrometry using a yeast test panel. *J CLIN MICROBIOL* 2014, 52, 3023–3029.
- [69] Wunschel, S.C., Jarman, K.H., Petersen, C.E., Valentine, N.B., et al., Bacterial analysis by MALDI-TOF mass spectrometry: an inter-laboratory comparison. *J. Am. Soc. Mass Spectrom.* 2005, 16, 456–62.
- [70] Mellmann, A., Bimet, F., Bizet, C., Borovskaya, A., et al., High interlaboratory reproducibility of matrix-assisted laser desorption ionization-time of flight mass spectrometry-based species identification of nonfermenting bacteria. *J CLIN MICROBIOL* 2009, 47, 3732–3734.
- [71] Chavy, A., Nabet, C., Normand, A.C., Kocher, A., et al., Identification of French Guiana sand flies using MALDI-TOF mass spectrometry with a new mass spectra library. *PLOS NEGLECT TROP D* 2019, 13, e0007031.
- [72] Li, Y., Wang, H., Zhao, Y.-P., Xu, Y.-C., Hsueh, P.-R., Evaluation of the Bruker biotyper matrix-assisted laser desorption/ionization time-of-flight mass spectrometry system for identification of *Aspergillus* species directly from growth on solid agar media. *FRONT MICROBIOL* 2017, 8, 1209.

- Application of proteome fingerprinting for fish stock monitoring
- 97.5% identification success
- Adjustment of identification scores improves identification success
- Line of evidence from different methods improves identification confidence
- Automatic mass spectra quality control

Journal Pre-proof