

Research



Cite this article: Juhel J-B *et al.* 2020

Accumulation curves of environmental DNA sequences predict coastal fish diversity in the coral triangle. *Proc. R. Soc. B* **287**: 20200248. <http://dx.doi.org/10.1098/rspb.2020.0248>

Received: 6 February 2020

Accepted: 18 June 2020

Subject Category:

Ecology

Subject Areas:

ecology

Keywords:

eDNA metabarcoding, sequence clustering, Operational Taxonomic Unit, diversity assessment, detectability

Author for correspondence:

Jean-Baptiste Juhel

e-mail: jeanbaptiste.juhel@gmail.com

†Joint last authorship.

Electronic supplementary material is available online at <https://doi.org/10.6084/m9.figshare.c.5047669>.

Accumulation curves of environmental DNA sequences predict coastal fish diversity in the coral triangle

Jean-Baptiste Juhel¹, Rizkie S. Utama², Virginie Marques¹, Indra B. Vimono², Hagi Yulia Sugeha², Kadarusman³, Laurent Pouyaud⁴, Tony Dejean⁵, David Mouillot^{1,6,†} and Régis Hocdé^{1,†}

¹MARBEQ, Univ. Montpellier, CNRS, Ifremer, IRD, Montpellier, France

²Research Center for Oceanography, Indonesian Institute of Sciences, Jl. Pasir Putih 1, Ancol Timur, Jakarta Utara, Indonesia

³Politeknik Kelautan dan Perikanan Sorong, KKD BP Sumberdaya Genetik, Konservasi dan Domestikasi, Papua Barat 98411, Indonesia

⁴Institut des Sciences de l'Évolution de Montpellier, Montpellier, France

⁵SPYGEN, 73370 Le Bourget-du-Lac, France

⁶ARC Centre of Excellence for Coral Reef Studies, James Cook University, Townsville, Australia

J-BJ, 0000-0003-2627-394X; VM, 0000-0002-5142-4191; K, 0000-0003-2312-2417; LP, 0000-0003-4415-9198; DM, 0000-0003-0402-2605; RH, 0000-0002-5794-2598

Environmental DNA (eDNA) has the potential to provide more comprehensive biodiversity assessments, particularly for vertebrates in species-rich regions. However, this method requires the completeness of a reference database (i.e. a list of DNA sequences attached to each species), which is not currently achieved for many taxa and ecosystems. As an alternative, a range of operational taxonomic units (OTUs) can be extracted from eDNA metabarcoding. However, the extent to which the diversity of OTUs provided by a limited eDNA sampling effort can predict regional species diversity is unknown. Here, by modelling OTU accumulation curves of eDNA seawater samples across the Coral Triangle, we obtained an asymptote reaching 1531 fish OTUs, while 1611 fish species are recorded in the region. We also accurately predict ($R^2 = 0.92$) the distribution of species richness among fish families from OTU-based asymptotes. Thus, the multi-model framework of OTU accumulation curves extends the use of eDNA metabarcoding in ecology, biogeography and conservation.

1. Introduction

Providing accurate biodiversity assessments is a critical goal in ecology and biogeography with estimations being constantly revised for some species-rich groups [1]. This issue is increasingly important, given the accelerating human footprint on Earth. The ongoing worldwide defaunation, characterized by massive population declines, may trigger the local or even global extinction of rare, elusive and cryptic species that are still unknown or poorly documented [2,3]. Such biodiversity losses directly impact ecosystem functioning, but also human health, well-being and livelihood [4,5]. This urges scientists to improve the accuracy and extend the breadth of biodiversity inventories and monitoring.

In the marine realm, the detection of species occurrences is particularly challenging due to the vast volume to monitor, the high diversity of habitats, the inaccessibility of some areas (e.g. deep sea) and the behaviour of some species (cryptobenthic or elusive) [6,7]. Environmental DNA (eDNA) metabarcoding is an emerging tool that can provide more accurate and wider biodiversity assessments than classical census methods, particularly for rare and elusive species [8–10]. This non-invasive method is based on retrieving DNA naturally released by organisms in their environment, amplified by polymerase chain reaction (PCR) and then sequenced to ultimately identify corresponding species [11]. However, inventorying and monitoring biodiversity using eDNA metabarcoding

requires the completeness of a reference database to accurately assign each sequence to a given species (e.g. [9]).

By now, only a minority of fish species are present in online DNA databases for mitochondrial regions targeted by metabarcoding markers, limiting the extent to which species diversity can be revealed by eDNA. This proportion of sequenced species is even lower in species-rich regions and poorly sampled habitats or taxa, while the effort to complete genetic reference databases is long and costly. As an alternative, a range of operational taxonomic units (OTUs) can be extracted from eDNA metabarcoding through filtering and clustering techniques [12]. Even if environmental genomics approaches have a long tradition of using OTU-based bioindicators [13], the extent to which the diversity of OTUs from a limited number of eDNA samples can reveal or predict the diversity of vertebrate species in a given biodiversity hotspot has not yet been investigated. This is particularly challenging for cryptobenthic fish species, which are key for reef ecosystems [14] but usually missed by classical surveys [7]. We thus urgently need a regional case study with a wide breadth of fish families and traits to test the potential of OTU-based assessment of biodiversity.

The Bird's Head Peninsula of West Papua (eastern Indonesia) is located in the centre of the Coral Triangle, which is known to host the world's richest marine biodiversity [15,16]. The current checklist of coastal fishes in the Bird's Head Peninsula identifies 1611 species belonging to 508 genera and 112 families [15,17], among which some are still poorly described or under severe threat [18–20]. Providing a blind but accurate assessment of the level and composition of a well-known vertebrate diversity from eDNA OTUs is thus a critical step in conservation, biogeography and ecology, particularly in such biodiversity hotspots.

Here, using eDNA metabarcoding from 92 seawater samples across the Bird's Head Peninsula, we (i) assessed the diversity of coastal fish species based on an online reference database for the teleo primers region of the 12S mitochondrial rDNA gene [21], (ii) estimated the diversity of fish OTUs based on a custom filtering and clustering bioinformatic pipeline, and (iii) tested the capacity of OTU accumulation curves to predict the level and composition of regional fish diversity.

2. Methods

(a) Sampling area and protocol

A total of 92 water samples were collected during October and November 2017 along the south coast of the Bird's Head region of West Papua (500 km) across different habitats but mainly coral reefs (electronic supplementary material, figure S1). Samples were collected in DNA-free plastic bags at the surface from a dinghy boat, at depths between 10 and 100 m during close circuit rebreather dives and at depths between 100 and 300 m using Niskin water samplers. A pressure and temperature sensor was coupled to the Niskin bottle to control the sampling depth and characterize the water mass via the vertical temperature profile. For each sample, 2 l of seawater was filtered with sterile Sterivex filter capsules (Merck Millipore; pore size 0.22 µm) using disposable sterile syringes. Immediately after, the filter units were filled with lysis conservation buffer (CL1 buffer SPYGEN) and stored in 50 ml screw-cap tubes at –20°C. A contamination control protocol was followed in both field and laboratory stages [21,22]. Water sample processing included the use of disposable gloves and single-use filtration equipment, and the bleaching (50% bleach) of Niskin water sampler.

(b) DNA extraction, amplification and high-throughput sequencing

The DNA extraction and amplification were performed following the protocol of [23], including 12 separate PCR amplifications per sample (see electronic supplementary material for more details on the protocol). A teleo-specific 12S mitochondrial rDNA primer (teleo, forward primer-ACACCGCCCGTCACTCT, reverse primer-CTTCCGGTACTTACCATG [21]) was used for the amplification of metabarcoding sequences, generating 63 ± 3 pb (mean \pm s.d.) long amplicons for all fish species referenced in EMBL database (European Molecular Biology Laboratory, www.ebi.ac.uk, v. 138, downloaded on January 2019) [24]. Eight negative extraction controls and two negative PCR controls (ultrapure water) were amplified (with 12 replicates as well) and sequenced in parallel to the samples to monitor possible contaminations. The teleo primers were 5'-labelled with an eight-nucleotide tag unique to each PCR replicate with at least three differences between any pair of tags, allowing the assignment of each sequence to the corresponding sample during sequence analysis. The tags for the forward and reverse primers were identical for each PCR replicate.

The purified PCR products were pooled in equal volumes, to achieve a theoretical sequencing depth of 1 000 000 reads per sample. Library preparation and sequencing were performed at Fasteris (Geneva, Switzerland). A total of five libraries were prepared using the MetaFast protocol (Fasteris, <https://www.fasteris.com/dna/?q=content/metafast-protocol-amplicon-metagenomic-analysis>), a ligation-based PCR-free library preparation. A paired-end sequencing (2 × 125 bp) was carried out using an Illumina HiSeq 2500 sequencer on three HiSeq Rapid Flow Cell v. 2 using the HiSeq Rapid SBS Kit v. 2 (Illumina, San Diego, CA, USA) following the manufacturer's instructions.

(c) Sequence analyses and taxonomic assignment

To evaluate the current completeness of the online database for the teleo region of the 12S mitochondrial DNA, an *in silico* PCR with 3 allowed mismatches using the teleo primers sequences was performed with ecoPCR [25] on the EMBL database. The generated list of sequenced species was compared with the checklists of fish species present in the Bird's Head region of Papua, provided by courtesy of Kulbicki *et al.* [17].

The amplified DNA sequences from the water samples were processed following two metabarcoding workflows. The first workflow used the OBITools software package [26] based on direct taxonomic assignment of the sequences using the ecotag lower common ancestor algorithm in EMBL database as a reference (see details in electronic supplementary material).

The ecotag algorithm can sometimes wrongly assign sequences to a given species or genus, despite a low-similarity percentage due to the incompleteness of reference database. We thus set the following similarity thresholds, 100–98, 90–98, 85–90 and 80–85% bp to assign sequences at the species, genus, family and order level, respectively. All the assignments with a similarity percentage lower than 80% were discarded from the analyses.

We evaluated the database completeness for the marker by running an *in silico* PCR on all fish mitochondrial DNA present in EMBL online database. A total of 394 species are sequenced in the Bird's Head region (24.5%, electronic supplementary material, table S1).

The second metabarcoding workflow was based on the SWARM clustering algorithm that groups multiple variants of sequences into OTUs [12]. Then, a post-clustering curation algorithm (LULU) was performed to curate data (see details in electronic supplementary material).

The SWARM clustering workflow was used to investigate the taxa present in the samples but not revealed by the taxonomic assignment process because of gaps in the EMBL database. The

number of taxa assigned in each family was corrected to avoid taxonomical redundancy assignment. For instance, the combined assignments to the genus *Zanclus* and the species *Zanclus cornutus* were considered as one taxa as potential PCR error may have produced two different assignment levels from the same sequence. These corrected numbers of taxa were then compared to the number of OTUs from the SWARM workflow in each family to evaluate the magnitude of the diversity missed by the direct assignment method. In the SWARM workflow, a family-level assignment was performed as well to remove the taxa that were not fish from non-specific amplifications and investigate the intrafamily diversity.

(d) Statistical analyses

To evaluate the number of taxa/OTUs present in the study area, a multi-model approach was implemented to fit asymptotes on the species and OTU accumulation curves. This approach considered five different accumulation models (Lomolino, Michaelis–Menten, Gompertz, asymptotic regression and logistic curve) and weighted them using the Akaike information criterion (AIC) [29]. For each curve, the accumulation model with the lowest AIC was selected. Accumulation curves and associated asymptotes were generated using the vegan R package. To estimate the sampling effort required to achieve a given proportion of asymptotes, we considered the model selected for accumulation curves. Then, we extracted the predicted number of samples producing a number of taxa/OTUs that outreached 90% and 95% of the asymptotes.

3. Results

(a) High heterogeneity of fish species detection among families

A total of 299 479 007 reads were produced using the OBITools pipeline over the 92 eDNA samples corresponding to 14 423 unique sequences with a mean of 307 unique sequences per sample (± 134 s.d.). In a conservative approach, stringent bioinformatic filters retained 9345 unique sequences, so 65% of the total. These 9345 unique sequences were then assigned to different taxonomic levels using the following genetic similarity thresholds: 98–100% for species, 90–98% for genus, 85–90% for family and 80–85% for order. This set of thresholds retained 7389 unique sequences resulting in 678 taxonomic assignments (electronic supplementary material, table S2).

A total of 310 species were detected, including 211 coastal fish species present in the checklist of the Bird's Head Peninsula and 99 fish species present in other regions but absent from this checklist (figure 1a). Conversely, 183 sequenced fish species which are present in the Bird's Head Peninsula were not detected in our eDNA samples using our stringent filters, representing 53.6% of the sequenced species present in the checklist. Since 75.5% of fish species in the checklist of the Bird's Head Peninsula were not sequenced for the 12S rDNA, the largest part of fish species diversity remained hidden through direct assignment (electronic supplementary material, table S1).

A total of 282 genera and 128 families of fish were detected compared with the regional checklist of 508 genera and 112 families out of which 46.1% and 72.3% are sequenced, respectively (electronic supplementary material, table S1). The number of fish species per family varied from 1 to 191 in the Bird's Head checklist (figure 1b), the richest family being the Gobiidae. Only 12 species of Gobiidae were detected in our 92 samples. Meanwhile, the most

represented family in the eDNA samples was the Labridae with 48 species (15.5% of the species found in the samples) out of 136 in the checklist (figure 1b).

The percentage of fish species sequenced per family varied between 0 and 100% with a mean of 40.3% ($\pm 31\%$ s.d.) in the Bird's Head Peninsula checklist while the percentage of detected species per family varied between 0 and 100% with a mean of 27.1% ($\pm 30.2\%$ s.d.) in eDNA samples (figure 1b). These two percentages were significantly and strongly related ($p < 0.001$) with the percentage of species sequenced per family explaining 85% of variation in the percentage of detected species per family (figure 1c).

(b) High but underestimated diversity of operational taxonomic units

Given that the low percentage of fish species sequenced for the 12S in the region is the main limitation to detect taxonomic diversity (figure 1c), we used an alternative approach based on unique clusters of genetic sequences called OTUs.

From the 331 839 591 initial reads, 4012 OTUs were generated using the SWARM clustering algorithm. After a series of post-clustering curation processes, 972 fish OTUs were filtered among which 819 were assigned to a family (electronic supplementary material, table S3). The number of detected OTUs varied from 1 to 54 among fish families (figure 2a), the richest families (greater than 50 OTUs) being the Gobiidae, Labridae and Pomacentridae. Overall, the number of OTUs was superior to the number of assigned taxa (genus and species) in 64.7% of the families found in the samples (mean $\Delta = 4 \pm 6.7$ s.d., figure 2a). This richness difference was null in 31.4% of the families and negative in 3.9% of them (figure 2a). This difference was notably high in some rich families such as the Gobiidae and Pomacentridae where the number of OTUs was more than 2 times and 1.5 times higher than the number of assigned taxa, respectively. By contrast, only 7 OTUs were produced compared with 11 assigned taxa for the Scombridae so $\Delta = -4$ units or -66.7% of this family richness.

The discrepancy between the two approaches (taxa and OTUs) was not significantly explained neither by the species richness of the family in the checklist ($R^2 < 0.01$, $p = 0.08$, figure 2b) nor by the percentage of sequenced fish species within each family in the checklist ($R^2 = 0.09$, $p = 0.05$, figure 2c).

On average, the number of OTUs underestimated the total number of coastal fish species in the Bird's Head Peninsula checklist with a mean net difference of 40.2% per family ($\pm 38.8\%$ s.d., figure 2d). For most families, this difference was high, reaching the maximum value of 95% for the Pseudochromidae. However, this difference could also be negative with more OTUs detected than species present in the checklist as for the Dasyatidae, Leiognathidae and Orectolobidae for which this difference reached -50% . Overall, the difference was marginally but significantly explained by the species richness of the family in the regional checklist ($R^2 = 0.09$, $p = 0.04$, figure 2d), suggesting that the bias is not proportional to the species richness of the family with species-rich families being more underestimated by OTUs than species-poor families.

(c) Prediction of fish species diversity from operational taxonomic unit accumulation curves

Since the two approaches (taxa and OTUs) underestimated the level of taxonomic diversity within fish families with a high

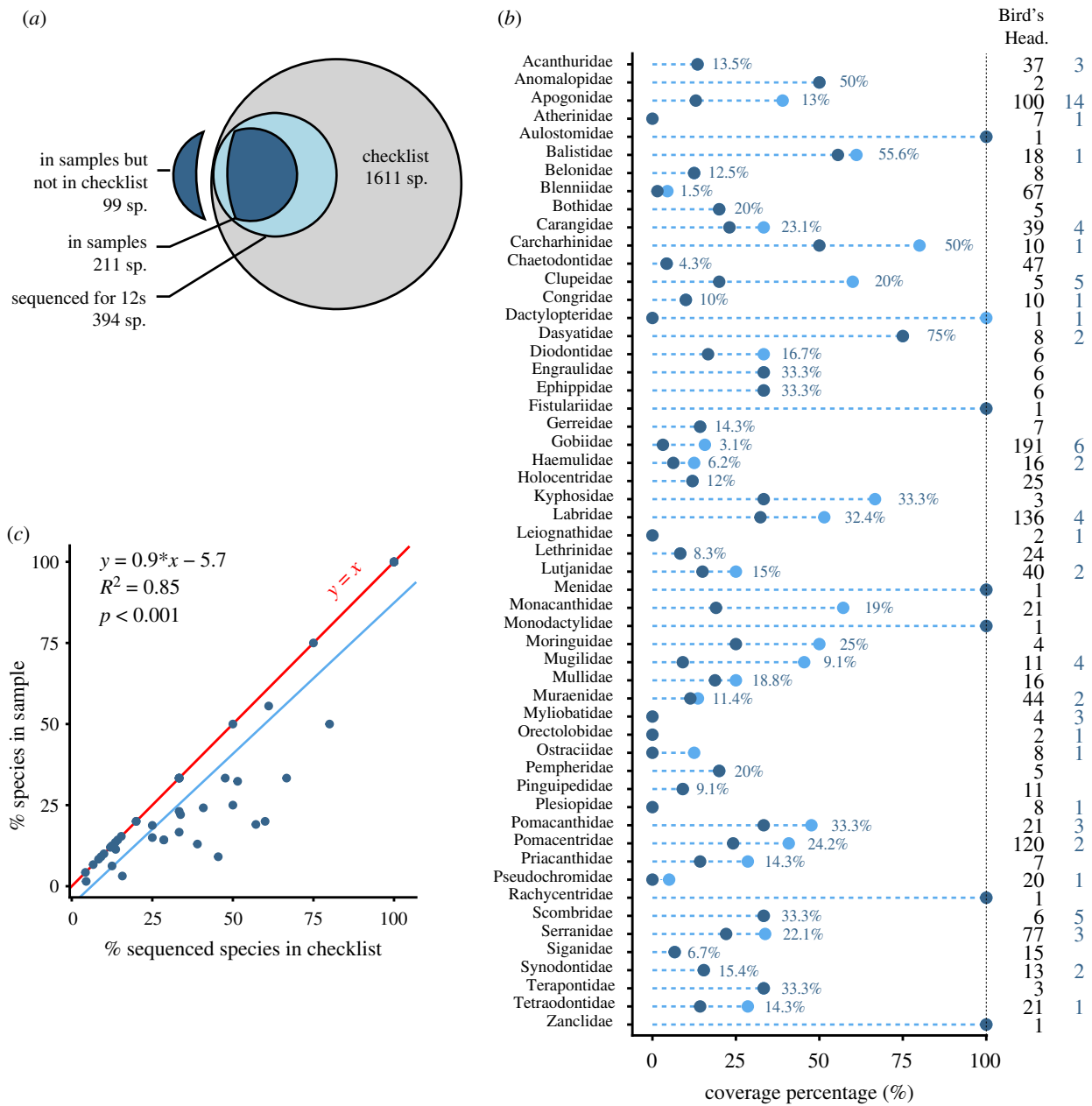


Figure 1. Number of fish species present in the checklist of the Bird's Head region (grey), sequenced in the European Molecular Biology Laboratory database (EMBL) (light blue) and detected in the eDNA samples (dark blue) (a); percentage of species detected in the samples (dark blue), sequenced in EMBL (light blue) in each family of species (b); percentage of species detected in the samples as a function of the percentage of sequenced species in EMBL (c). (b) The percentages of the species detected in the eDNA samples compared with the species present in the Bird's Head region are displayed next to the points. The number of species per family in the checklist and the number of species detected in the samples but not present in the checklist are both on the right of the figure in black and dark blue, respectively. Only the sequences assigned to species using ecotag program (similarity >98%) are used in this figure. (c) Each point corresponds to a fish family. (Online version in colour.)

uncertainty, we modelled accumulation curves from the diversity of species and OTUs found across our 92 samples. The modelled asymptote of the assigned species reached 429 species, a value very close to the 394 sequenced species present in the Bird's Head Peninsula, but 3.7 times lower than the 1611 species in the regional checklist (figure 3a). Meanwhile, the OTU accumulation curve reached an asymptote of 1531; a value close (95%) to the number of fish species (1611) referenced in the checklist of the Bird's Head Peninsula.

Applying this method to the 15 fish families which counted more than 10 OTUs and 10 species in the checklist permitted to assess the ability of eDNA-based accumulation curves to predict regional fish richness. For instance, the

OTU accumulation curves for the Gobiidae, Labridae and Pomacentridae, the three richest families (51, 54 and 53 OTUs, respectively), produced asymptotes and thus predictions of fish diversity much lower than those in the regional checklists with 107.5, 66.1 and 76.2 OTUs (i.e. 47.5%, 81.7% and 69.6% of the checklist richness respectively; figure 3b-d).

We then tested the ability of the assigned taxa, the OTUs and the OTU accumulation curve approaches to predict fish species richness within families of the regional checklist, so the predictive power of linear or proportional relationships. The total number of assigned taxa per family in our samples was a significant but weak predictor of the number of fish species per family in the checklist ($R^2 = 0.60$, $p < 0.001$,

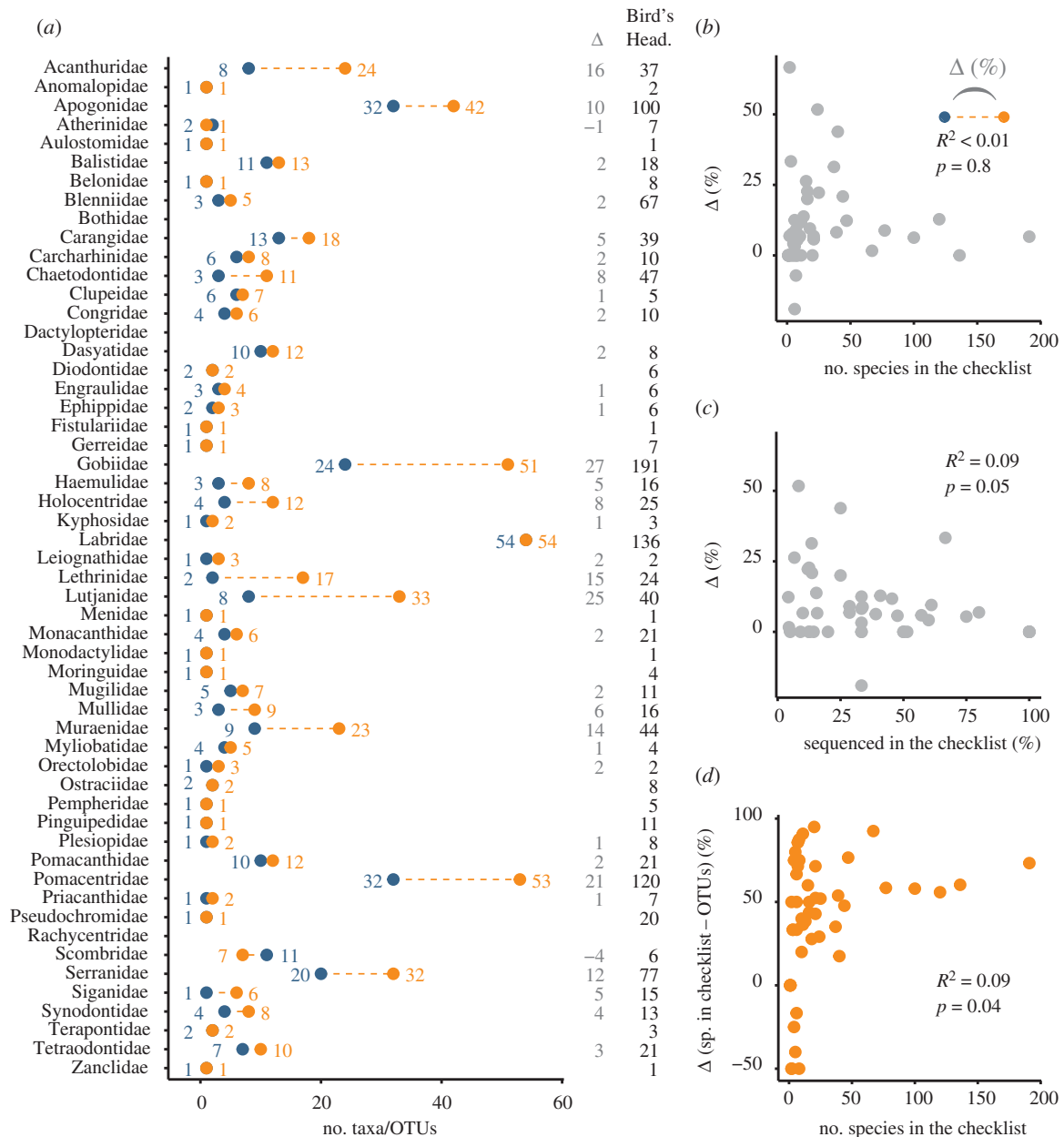


Figure 2. Number of taxa assigned by the OBITools workflow (blue) and number of OTUs generated by the SWARM workflow (orange) in the different fish families (a); distribution of the differences between the two workflows as a function of family richness (b) and as a function of family sequencing coverage (c); distribution of the differences between OTUs and the number of species in the checklist as a function of family richness (d). (a) The difference of taxa/OTUs between the two methods (noted Δ) and the number of species in the checklist of the Bird's Head region are on the right of the figure in grey and black, respectively. For the OBITools workflow, only the sequences assigned to species and genus using ecotag program (similarity greater than 98% and greater than 90% respectively) are used in this figure. For the SWARM workflow, only the OTUs curated by LULU and assigned to family (similarity greater than 85%) are used in this figure. (Online version in colour.)

figure 4a) with the richness of some families being largely underestimated (e.g. 87.4% of net difference with the checklist for the Gobiidae, figure 4a,d). The number of OTUs per family was a better predictor of the family species richness in the checklist ($R^2 = 0.80$, $p < 0.001$) but left 20% of unexplained variation among families with still a marked underestimation (73.3% of net difference with the checklist for Gobiidae, figure 4b,e). Using the asymptotes of OTU accumulation curves, we obtained a high predictive accuracy of $R^2 = 0.92$ ($p < 0.001$) for the species richness within families with less bias for the Gobiidae (43.7% of net difference with the checklist) (figure 4c,f).

In addition, we observed that the net difference between the number of assigned taxa per family and the number of

species per fish family in the checklist is not related to the number of species of the families (figure 4d), suggesting an absence of systematic bias towards the underestimation of species-rich families. By contrast, the net difference between the number of OTUs per fish family and the number of species per family in the checklist significantly increased ($R^2 = 0.35$, $p = 0.02$) with the number of species per family (figure 4e). This bias towards the underestimation of species richness within species-rich families is nonetheless avoided when using the asymptotes of OTU accumulation curves ($p = 0.24$, figure 4f). Thus, asymptotes of OTU accumulation curves are most accurate and least biased eDNA-based predictors of fish species diversity within families in this marine biodiversity hotspot.

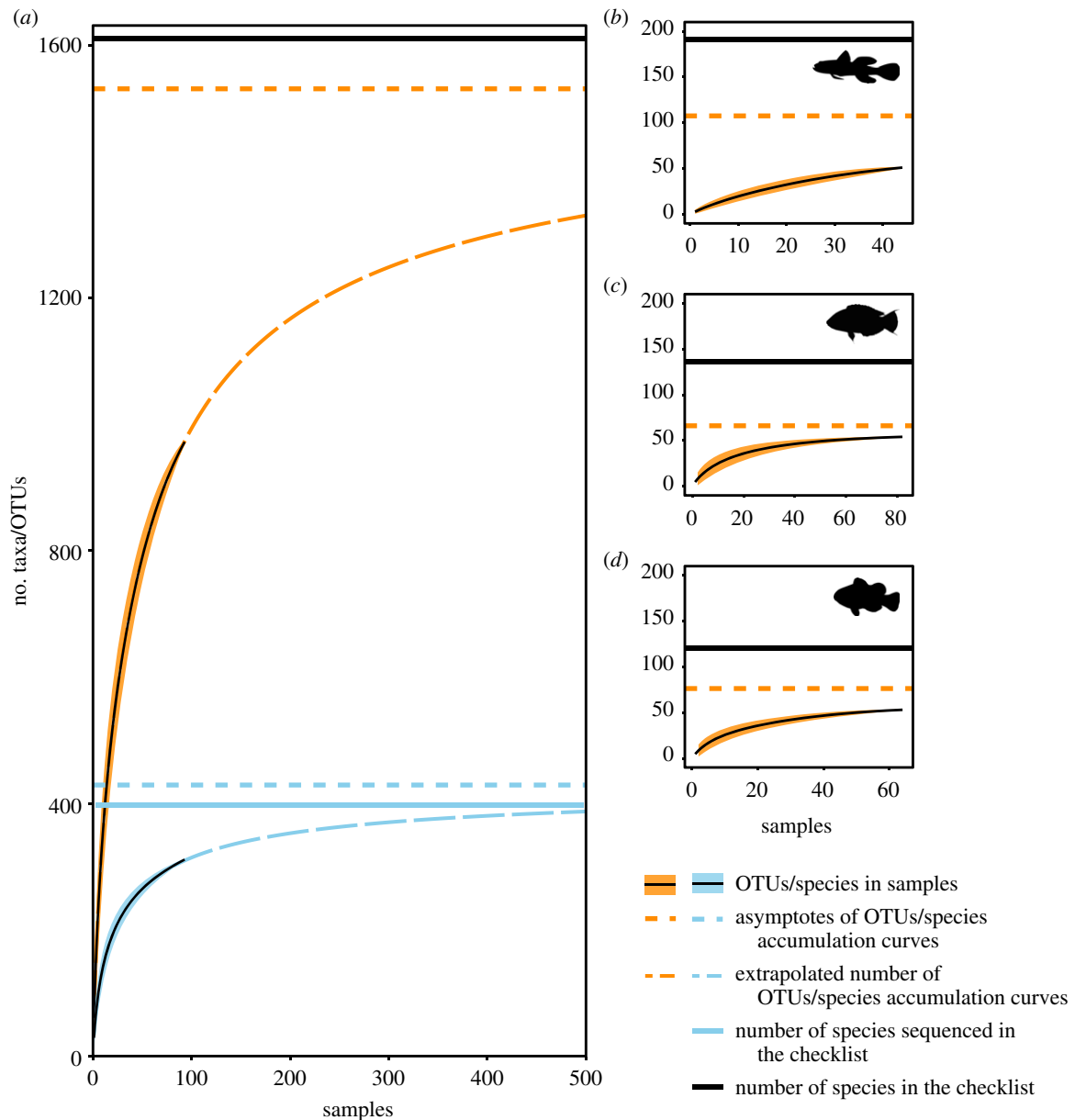


Figure 3. Accumulation curves of species assigned (blue) and the OTUs (orange) obtained in the whole sampling (a) and within the three most diverse families: Gobiidae (b), Labridae (c) and Pomacentridae (d). The detection of species and OTUs was randomized 100 times and the results were used to generate the confidence intervals. The asymptotes were modelled by a multi-model approach weighted by the Akaike information criterion (AIC). Fish silhouettes are from phylopic.org (Kent Sorgon & Lily Hughes). (Online version in colour.)

(d) Sampling efforts necessary to achieve regional fish diversity inventory

Not only the OTU accumulation curves and their asymptotes provide diversity estimates, they also provide crucial insights into the sampling effort needed to achieve a more complete census. Here, using the asymptote on the OTU accumulation curve for all fish species (figure 3a), we found that our 92 cumulated samples (representing 0.2 m³) achieved up to 63.5% of the potential fish OTU diversity in the Bird's Head Peninsula (figure 5). To collect 90% of this regional fish diversity, we should have filtered seawater in 735 samples, so eight times the effort of our sampling campaign, representing an aggregated sampled water volume of 1.5 m³. This sampling effort would reach 1883 samples (an aggregated water volume of 3.8 m³) to collect 95% of the regional fish OTU richness (figure 5).

On average across fish families, our sampling effort achieved the detection of 77.1% (± 14.9 s.d.) of OTUs predicted by the asymptote of the accumulation curve with a variation

among families ranging from 42.2% (Muraenidae) and 47.5% (Gobiidae) to 93.9% (Balistidae) (figure 5). The sampling effort needed to achieve 90% of the asymptotic number of OTUs in the region varied greatly among families, ranging from 37 samples for Chaetodontidae to 494 samples for Gobiidae, with a mean of 164 samples (± 123 s.d.). The estimated additional sampling effort to reach 95% from 90% of the OTU richness ranged from 20 more samples (Tetraodontidae) to 593 more samples (Gobiidae).

4. Discussion

(a) Overcoming incompleteness of genetic reference databases

Environmental DNA metabarcoding has the potential to surpass most classical survey methods to assess biodiversity in both terrestrial and aquatic systems [30]. Yet, genetic reference

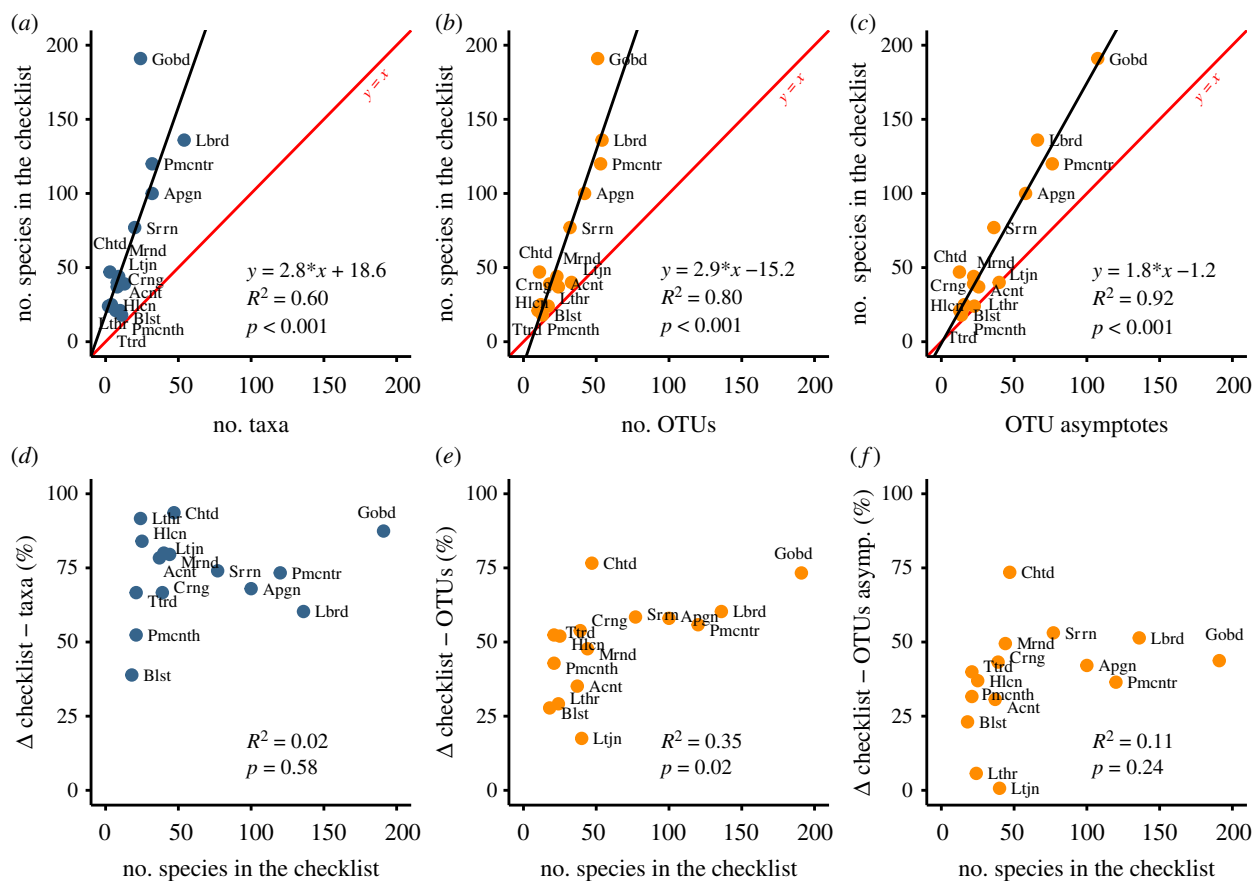


Figure 4. Linear regression of the diversity of the most diverse families as a function of the number taxa assigned (a), the number of OTU (b), the asymptotes of the OTU accumulation curves (c) and differences between the number of taxa assigned (d), the number of OTUs (e), the asymptotes of OTU accumulation curves (f) and the number of species in the checklist as a function of the number of species in the checklist. Only the families with a number of OTU and a number of species in the checklist greater than or equal to 10 are presented to provide accurate estimations. (Online version in colour.)

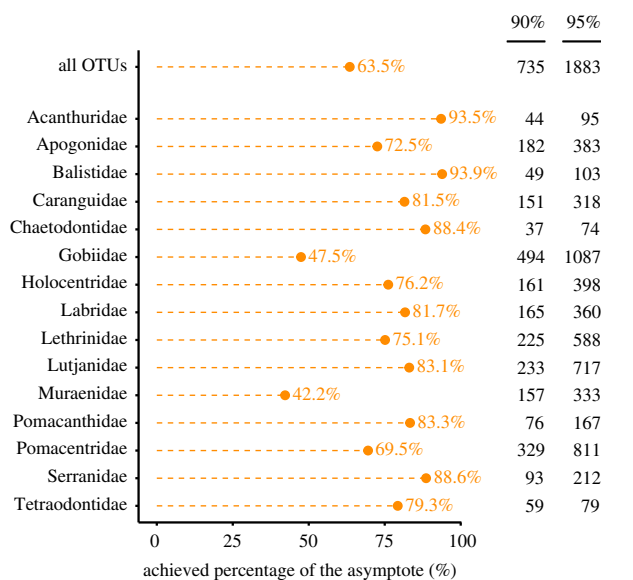


Figure 5. Percentage of the OTUs diversity covered by the current sampling effort ($n = 92$) in the families of fish (orange) and the estimated sampling effort required to achieve both 90% and 95% of the diversity. Only the families with a number of OTU and a number of species in the checklist greater than or equal to 10 are presented to provide accurate estimations. (Online version in colour.)

databases are often incomplete, especially for species-rich ecosystems such as the Coral Triangle, a global marine biodiversity hotspot [14]. For instance, the current completeness

of the 12S rDNA online databases for the teleo primer covers only 24.5% of fish species in the Bird's Head Peninsula. Meanwhile, this cover reaches 77.3% for the COI (mitochondrial cytochrome *c* oxidase subunit I), but fish COI primers still perform poorly in comparison to 12S markers [31].

With around 28% of families, 54% of the genera and 76% of species not sequenced for the 12S rDNA teleo primers region, the largest part of fish diversity in the Bird's Head Peninsula remains thus hidden through direct assignment. Additionally, sequences present in the reference online databases may have been collected from individuals not located in the region of interest. This can induce assignment errors due to biogeographic-related genetic variation (e.g. [32]). The lack of sequencing coverage highlights the immense gap to be filled for online databases to be exhaustive, while numerous species still remain to be described [33]. This limitation prevents metabarcoding approaches from characterizing entire fish assemblages through direct species assignment. Yet, the tax-assignment method reveals the presence of 211 fish species referenced in the checklist of coastal fishes in the Bird's Head Peninsula (figure 1a). Conversely, 99 assigned species were absent from this checklist. These 99 detections can either be true presences extending the distribution of some species and revisiting the regional checklist or false presences due to wrong assignments or possible contaminations. For instance, the Atlantic salmon (*Salmo salar*), probably a laboratory kit contaminant, was found in our study and removed from the analyses (see Methods). The large number of species present

in the samples but absent from the regional checklist suggests that inventories of some families are still incomplete. On average, 2.5 detected species per family (± 2.6 s.d., figure 1b) are absent from the checklist, ranging from 0 to 14 species (Apogonidae). This mismatch allows us to target future sampling efforts towards families and their habitats to complete the regional checklist.

As an alternative to species assignment, the use of OTUs as species proxy units is an option that has not yet been tested for vertebrates in species-rich ecosystems while currently used when the concept of species is debatable like for fungi or unicellular organisms [34,35].

Here, using a conservative and stringent bioinformatic pipeline, we show that the diversity of OTUs is a weak and biased estimator of species diversity with species-rich families being strongly underrepresented. To overcome this limitation, we propose to rely on OTU accumulation curves which provide an unbiased estimate of regional fish diversity and fish richness within families. The asymptotes underestimate the regional fish species richness, but the bias is highly consistent among families (figure 4f). We thus propose to extend this method for taxonomic inventories in poorly sampled ecosystems like the deep sea to estimate the diversity at different taxonomic levels.

(b) Revealing the potential and limitation of eDNA metabarcoding inventories

Fishes are the most diverse group of vertebrates on Earth with varying body sizes, environmental niches and diets. Monitoring fish assemblages in marine biodiversity hotspots like the Coral Triangle is a great challenge, particularly for small, rare, cryptobenthic or elusive species. Here, we show that the percentage of sequenced species is highly variable among families preventing any robust estimation of species richness. Instead, OTUs have the potential to reveal the presence of a broad range of fish species (i.e. from different lineages and with contrasted life-history traits). For instance, cryptobenthic families have been poorly documented and are often ignored in traditional visual censuses [7], while they strongly influence ecosystem functioning [13]. Similarly, traditional visual censuses often miss highly mobile and elusive species such as sharks [9].

Among the 310 assigned fish species, we detected the presence of small cryptobenthic species such as *Gobiodon histrio* or *Ostorhinchus selas*, a goby and a cardinalfish with a maximum length below 40 mm, respectively. We also detected large pelagic fish such as the dogtooth tuna (*Gymnosarda unicolor*) or the thresher shark (*Alopias pelagicus*) reaching over 2 m and 4 m long, respectively. Flagship species for conservation were also present in our DNA samples such as the over-exploited Napoleon wrasse (*Cheilinus undulatus*, Endangered, IUCN Red List, www.iucnredlist.org), the Scalloped hammerhead shark (*Sphyrna lewini*, Endangered) and several shark species being classified as Near Threatened (NT) (*C. brevipinna*, *C. Leucas*, *C. sorrah*, *C. melanopterus*, *T. obesus*).

Even if not assigned at species level, OTUs can be defined as distinct entities for which their distribution and temporal variability can be assessed and monitored [36]. Moreover, the OTUs and their associated sequences can remain in public repositories until they are assigned to a species, subspecies or complex as databases improve [37]. However, the major caveat of using OTUs for diversity inventories is that

they cannot be directly considered as species with complete certainty. Species with intra-specific genetic variability can produce two separate OTUs, overestimating species diversity. Conversely, two species phylogenetically close to each other with low genetic variability can be grouped into a single OTU, thus underestimating species diversity. The accuracy of diversity inventories using eDNA metabarcoding is thus directly based on the taxonomic resolution of the barcode used and genetic variability among families but also the number of samples.

Here, we also reveal the gap of biodiversity that remains to be detected using OTU accumulation curves. The effort can be massive for some families (figure 5) and more ambitious eDNA sampling campaigns should be on the agenda in species-rich regions like the Coral Triangle. OTU accumulation curves can also serve to evaluate the efficiency of a sampling method (e.g. punctual filtration, transect filtration), the sampled area or the diversity of habitats that are required (e.g. depth, complexity, distance from the seafloor) and their location (e.g. proximity of reefs, hotspots) especially when targeting rare, elusive, highly mobile or cryptobenthic families of fish.

The contrasts between assigned taxa diversity, OTU diversity and OTU asymptote diversity show that the detectability varies strongly among fish families. These contrasts can be related to the ecology of the species but also to the state of the retrieved DNA fragments (intra or extracellular), their sources (e.g. gametes, larvae, faeces), their release rate, their diffusion in the water column (limited or wide) and their transportation [38]. For instance, benthic fish species such as gobies with a small movement range would release DNA fragments through skin and faeces on a small area. However, such species could release a massive number of gametes carried through the water column [13] so may appear highly detectable during breeding season. Further comparative works are urgently needed between visual, camera and eDNA metabarcoding surveys to better estimate the level of detectability of each species or family in order to provide reliable biodiversity assessments. For instance, coupling eDNA metabarcoding and video surveillance allows the detection of 82 fish genera from 13 orders on reefs and seagrass with only 24 genera in common [39]. Investigating biodiversity should also consider its multiple components including functional and phylogenetic diversity that are key for reef ecosystem functioning [40]. Associating OTUs to species might allow us to fill this gap, but it will require massive sampling and sequencing efforts.

Data accessibility. The metadata and bioinformatic outputs are available in the Dryad Digital Repository [41]. The metabarcoding pipelines are available in GitLab (https://gitlab.mbb.univ-montp2.fr/edna/snakeyaml_only_obitools and https://gitlab.mbb.univ-montp2.fr/edna/bash_swarm).

Authors' contributions. J.-B.J., I.B.V., K., L.P., D.M. and R.H. designed research; J.-B.J. and R.H. design the specific research methods of data collection and the sampling strategy; J.-B.J., R.S.U., K. and R.H. collected samples and data; T.D. coordinated the biomolecular analyses; J.-B.J., R.S.U. and V.M. performed the bioinformatics analyses; J.-B.J., R.S.U., V.M., T.D., L.P., D.M. and R.H. defined sequencing strategy, analysed and interpreted data; J.-B.J. wrote the initial draft and designed the figures; J.-B.J., R.S.U., V.M., I.B.V., H.Y.S., K., T.D., L.P., D.M. and R.H. wrote the paper and approved the final draft; and L.P., D.M. and R.H. acquired funding to conduct the study.

Competing interests. We declare we have no competing interests.

Funding. Fieldwork and laboratory activities were supported by the Lengguru 2017 Project (www.lengguru.org), conducted by the French National Research Institute for Sustainable Development

(IRD), the Indonesian Institute of Sciences (LIPI) with the Research Center for Oceanography (RCO, the Politeknik KP Sorong), the University of Papua (UNIPA) with the help of the Institut Français en Indonesia (IFI) and with corporate sponsorship from the Total Foundation and TIPCO company. The eDNA sequencing was funded by Monaco Explorations.

Acknowledgements. We thank the Indonesian Institute of Sciences (LIPI) for promoting our collaboration and the Sorong Polytechnic of Marine and Fisheries (Politeknik KP Sorong, West Papua) for providing the vessel *Airaha 02* that we used in this campaign. We thank the crew of the *Aihara 02* for assisting us during the operations and the SPYGEN staff for the technical support in the laboratory.

References

- Costello MJ, Chaudhary C. 2017 Marine biodiversity, biogeography, deep-sea, and conservation. *Curr. Biol.* **27**, R511–R527. (doi:10.1016/j.cub.2017.04.060)
- Barlow J *et al.* 2018 The future of hyperdiverse tropical ecosystems. *Nature* **559**, 517–526. (doi:10.1038/s41586-018-0301-1)
- Lees AC, Pimm SL. 2015 Species, extinct before we know them. *Curr. Biol.* **5**, R177–R180. (doi:10.1016/j.cub.2014.12.017)
- Díaz S *et al.* 2018 Assessing nature's contributions to people. *Science* **359**, 270–272. (doi:10.1126/science.aap8826)
- Duffy JE, Godwyn CM, Cardinale BJ. 2017 Biodiversity effects in the wild are common and as strong as key drivers of productivity. *Nature* **549**, 261–264. (doi:10.1038/nature23886)
- Juhel JB, Vigliola L, Wantiez L, Letessier TB, Meeuwig JJ, Mouillot D. 2019 Isolation and no-entry marine reserves mitigate anthropogenic impacts on grey reef shark behavior. *Sci. Rep.* **9**, 2897. (doi:10.1038/s41598-018-37145-x)
- Brandl SJ, Goatley CHR, Bellwood DR, Tornabene L. 2018 The hidden half: ecology and evolution of cryptobenthic fishes on coral reefs. *Biol. Rev.* **93**, 1846–1873. (doi:10.1111/brv.124233)
- Garlapati D, Charankumar B, Ramu K, Madeswaran P, Ramana Murthy MV. 2019 A review on the applications and recent advances in environmental DNA (eDNA) metagenomics. *Rev. Environ. Sci. Bio.* **18**, 389. (doi:10.1007/s11157-019-09501-4)
- Boussarie G *et al.* 2018 Environmental DNA illuminates the dark diversity of sharks. *Sci. Adv.* **4**, eaap9661. (doi:10.1126/sciadv.aap9661)
- Fukamoto S, Ushimaru A, Minamoto T. 2015 A basin-scale application of environmental DNA assessment for rare endemic species and closely related exotic species in rivers: a case study of giant salamanders in Japan. *J. Appl. Ecol.* **52**, 358–365. (doi:10.1111/1365-2664.12392)
- Ruppert KM, Kline RJ, Rahman MDS. 2019 Past, present, and future of environmental DNA (eDNA) metabarcoding: a systematic review in methods, monitoring, and applications of global eDNA. *Glob. Ecol. Conserv.* **17**, e00547. (doi:10.1016/j.gecco.2019.e00547)
- Mahé F, Rognes T, Quince C, de Vargas C, Dunthorn M. 2014 Swarm: robust and fast clustering method for amplicon-based studies. *PeerJ* **2**, e593. (doi:10.7717/peerj.593)
- Cordier T *et al.* 2019 Multi-marker eDNA metabarcoding survey to assess the environmental impact of three offshore gas platforms in the North Adriatic Sea (Italy). *Mar. Environ. Res.* **146**, 24–34. (doi:10.1016/j.marenvres.2018.12.009)
- Brandl SJ, Rasher DB, Côté IM, Casey JM, Darling ES, Lefcheck JS, Duffy JE. 2019 Coral reef ecosystem functioning: eight core processes and the role of biodiversity. *Front. Ecol. Environ.* **17**, 445–454. (doi:10.1002/fee.2088)
- Veron JEN, Devantier LM, Turak E, Green AL, Kininmonth S, Stafford-Smith M, Peterson N. 2009 Delineating the coral triangle. *Galaxea, JCRS* **11**, 91–100. (doi:10.3755/galaxea.11.91)
- Allen GR, Erdmann MV. 2012 *Reef fishes of the East Indies*. Volumes I–III. Perth, Australia: Tropical Reef Research.
- Kulbicki M *et al.* 2013 Global biogeography of reef fishes: a hierarchical quantitative delineation of regions. *PLoS ONE* **8**, e81847.
- Exton DA *et al.* 2019 Artisanal fish fences pose broad and unexpected threats to the tropical coastal seascape. *Nat. Commun.* **10**, 2100. (doi:10.1038/s41467-019-10051-0)
- Jones LA, Mannion PD, Farnsworth A, Valdes PJ, Kelland S-J, Allison PA. 2019 Coupling of palaeontological and neontological reef coral data improves forecasts of biodiversity responses under climatic change. *R. Soc. Open Sci.* **6**, 182111. (doi:10.1098/rsos.182111)
- Ainsworth CH, Pitcher TJ, Rotinsulu, C. 2008 Evidence of fishery depletions and shifting cognitive baselines in Eastern Indonesia. *Biol. Conserv.* **141**, 848–859. (doi:10.1016/j.biocon.2008.01.006)
- Valentini A *et al.* 2016 Next-generation monitoring of aquatic biodiversity using environmental DNA metabarcoding. *Mol. Ecol.* **25**, 929–942. (doi:10.1111/mec.13428)
- Goldberg CS *et al.* 2016 Critical considerations for the application of environmental DNA methods to detect aquatic species. *Methods Ecol. Evol.* **7**, 1299–1307. (doi:10.1111/2041-210X.12595)
- Pont D *et al.* 2018 Environmental DNA reveals quantitative patterns of fish biodiversity in large rivers despite its downstream transportation. *Sci. Rep.* **8**, 10361. (doi:10.1038/s41598-018-28424-8)
- Baker W, van den Broek A, Camon E, Hingamp P, Sterk P, Stoesser G, Tuli MA. 2000 The EMBL nucleotide sequence database. *Nucleic Acids Res.* **28**, 19–23. (doi:10.1093/nar/gki098)
- Ficetola GT, Coissac E, Zundel S, Riaz T, Shehzad W, Bessi re J, Taberlet P, Pompanon F. 2010 An *in silico* approach for the evaluation of DNA barcodes. *BMC Genomics* **11**, 434. (doi:10.1186/1471-2164-11-434)
- Boyer F, Mercier C, Bonin A, Le Bras Y, Taberlet P, Coissac E. 2016 OBITOOLS: a UNIX-inspired software package for DNA metabarcoding. *Mol. Ecol. Res.* **16**, 176–182. (doi:10.1111/1755-0998.12428)
- Larkin MA *et al.* 2007 Clustal W and Clustal X version 2.0. *Bioinformatics* **23**, 2947–2948. (doi:10.1093/bioinformatics/btm404)
- Kearse M *et al.* 2012 Geneious basic: an integrated and extendable desktop software platform for the organization and analysis of sequence data. *Bioinformatics* **28**, 1647–1649. (doi:10.1093/bioinformatics/bts19)
- Aho K, Derryberry D, Peterson T. 2014 Model selection for ecologists: the worldviews of AIC and BIC. *Ecology* **95**, 631–636. (doi:10.1890/13-1452.1)
- Deiner K *et al.* 2017 Environmental DNA metabarcoding: transforming how we survey animal and plant communities. *Mol. Ecol.* **26**, 5872–5895. (doi:10.1111/mec.14350)
- Collins RA, Bakker J, Wangenstein OS, Soto AZ, Corrigan L, Sims DW, Genner MJ, Mariani S. 2019 Non-specific amplification compromises environmental DNA metabarcoding with COI. *Methods Ecol. Evol.* **10**, 1985–2001. (doi:10.1111/2041-210X.13276)
- Wadrop E, Hobbs J-P, Randall JE, DiBattista JD, Rocha LA, Kosaki RK, Berumen ML, Bowen BW. 2016 Phylogeography, population structure and evolution of coral-eating butterflyfishes (Family Chaetodontidae, genus *Chaetodon*, subgenus *Corallochaetodon*). *J. Biogeogr.* **43**, 1116–1129. (doi:10.1111/jbi.12680)
- Pinheiro HT, Moreau S, Daly M, Rocha LA. 2019 Will DNA barcoding meet taxonomic needs? *Science* **365**, 873–875. (doi:10.1126/science.aay7174)
- Pawlowski J *et al.* 2018 The future of biotic indices in the ecogenomic era: integrating (e)DNA metabarcoding in biological assessment of aquatic ecosystems. *Sci. Total Environ.* **637–638**, 1295–1310. (doi:10.1016/j.scitotenv.2018.05.002)
- Lladó FS, Větrovský T, Baldrian P. 2019 The concept of operational taxonomic units revisited: genomes of bacteria that are regarded as closely related are often highly dissimilar. *Folia Microbiol.* **64**, 19–23. (doi:10.1007/s12223-018-0627-y)
- Cordier T, Esling P, Lejzerowicz F, Visco J, Ouadahi A, Martins C, Cedhagen T, Pawlowski J. 2017 Predicting the ecological quality status of marine environments from eDNA metabarcoding data using supervised machine learning. *Environ. Sci. Technol.* **51**, 9118–9126. (doi:10.1021/acs.est.7b01518)

37. Wangenstein O, Palacín C, Guardiola M, Turon X. 2018 DNA metabarcoding of littoral hard-bottom communities: high diversity and database gaps revealed by two molecular markers. *PeerJ* **6**, e4705. (doi:10.7717/peerj.4705)
38. Harrison JB, Sunday JM, Rogers SM. 2019 Predicting the fate of eDNA in the environment and implications of studying biodiversity. *Proc. R. Soc. B* **286**, 20191409. (doi:10.1098/rspb.2019.1409)
39. Stat M, Jeffrey J, DiBattista JD, Newman SJ, Bunce M, Harvey ES. 2018 Combined use of eDNA metabarcoding and video surveillance for the assessment of fish biodiversity. *Conserv. Biol.* **33**, 196–205. (doi:10.1111/cobi.13183)
40. Duffy JE, Lelcheck JS, Stuart-Smith RD, Navarrete SA, Edgar GJ. 2016 Biodiversity enhances reef fish biomass and resistance to climate change. *Proc. Natl Acad. Sci. USA* **113**, 6230–6235. (doi:10.1073/pnas.1524465113)
41. Juhel J-B *et al.* 2020 Data from: Accumulation curves of environmental DNA sequences predict coastal fish diversity in the coral triangle. Dryad Digital Repository. (doi:10.5061/dryad.t1g1jw05)