

## **REVIEWER 1**

Dear Editor,

Please receive my report for the proposed paper, which is very positive given the importance and presentation style of the described data.

*We thank the reviewer for their constructive comments and overall endorsement of the work. Changes have been made to the wording in the text following suggestions by both reviewers. All line numbers in the response to reviewers refer to the resubmitted manuscript and all changes to the text are provided as track changes. Detailed responses to each of the reviewers' comments are provided below.*

My "minor comments" are indicated below:

Line 85. "spatial distributions were generated for each taxon from four biotic groups ". I suggest to clearly define the habitats for the groups with some sort of seascape approach in which the general features of the dominating geomorphology are considered.

*Different taxa within each biotic group can occur across a range of habitats or dominating geomorphologies. Given the number of taxa modelled, it is therefore difficult to provide a high level overview for each biotic group that specially links to seascapes / dominating geomorphology. However, we have strived to add some clarity on habitats associated with the various taxonomic groups where possible (L 85 - 88):*

*"Briefly, spatial distributions were generated for each taxon from four biotic groups: demersal fish (n taxa = 235, bottom dwelling fish occurring on, or near, the seafloor throughout the study area), reef fish (n taxa = 51, reef associated fish, occurring on coastal rocky reefs), subtidal invertebrates (n taxa = 207, bottom dwelling or benthic invertebrates occurring throughout the study area) and macroalgae (n taxa = 86, occurring in the coastal margin within the photic zone) using ensemble SDM methods."*

"Demersal fish", this name is referring to fishery resources. Therefore, I would suggest to disambiguate that diction by stating that authors are referring to muddy bottoms in continental margins. Also, "reef fish" is ambiguous since "reef habitat" (line 94) are nor clearly defined. There are rocky reefs but also biogenic reefs (e.g. deep-sea corals).

*We believe that demersal fish refers to fish occurring in the demersal zone (the water column near to the seabed) rather than as a fishery term exclusively. We now clarify this, and the 'reef habitat' at L 85 – 88 (in response to the previous reviewer comment).*

Line 85 and Figure 1. Please revise the N (235 vs 239 for Demersal fish taxa

*Thank you – this has been changed to 235 (L85)*

Lines 87-89. For all names and acronyms, please choose if place the first letter in majuscule or not throughout the MS.

*Thank you for the comment – we consistently capitalise some words throughout the document (e.g., Boosted Regression Tree). We leave it to the copy editor to decide whether to retain as capitals or not (we have seen this term written both capitalised or not).*

Line 90, about figure 2. In that figure:

1-To better appreciate the message I would provide a plate collating A and B into a zoomed sub-area.

2-I feel missing a plate where data collection geographic points are indicated in a third plate. In that manner the whole conceptual data elaboration/modeling process could be appreciated.

*We are not completely clear about the reviewers first suggestion, but if it relates to being able to zoom in to look at the detail of the predictions, we believe this is best done online from the New Zealand Department of Conservation GIS webviewer. We now emphasise this point in the figure caption which reads (L 143)*

*“Figure 2. Example output map from the atlas of seabed biodiversity for Aotearoa New Zealand as plotted in the NZ Department of Conservation’s online geoportal available at: <https://doc-marine-data-deptconservation.hub.arcgis.com/>. The example includes; (A) mean predicted habitat suitability index (HSI, 0 - 1) for Australasian snapper (*Chrysophrys auratus*, a demersal fish species) and (B) associated estimates of uncertainty (measured as the standard deviation of the mean HSI). The geoportal provides an easy means to explore the data, including by zooming into areas of interest (see supplementary materials 2). Image produced under permission from the NZ Department of Conservation.”*

*We appreciate the reviewers’ suggestions. We can produce a third panel with the data points for the example taxa in Figure 2 but we do not make the raw point data available as part of this work (most are easily accessible from the references provided in the supplementary materials but some of these data are proprietary – which we may be able to make available upon reasonable request). Therefore, we have decided to omit this additional figure panel as it may be misleading as to the available data. The dataset presented for publication is focussed on the spatial predictions all of which are freely accessible.*

Line 105. “20 spatially explicit environmental variables”. The characterisation of a habitat envelope is a fundamental importance for modelling species distributions in areas with no sampling data. More details should be provided in the supplementary material about environmental data extraction. For example, how authors derived Seabed Temperature and Salinity based on surface (SST) data? Authors quote in the table "Oceanographic data from CARS2009 (NIWA, unpublished 2009). This is quoted also for other environmental variables. More details would be welcome here on that procedure ("Derived from SST described above at two resolutions and merged"??). Also, it may be beneficial to have BotTemp and SSTGad sequentially together in the Supplementary Table.

*The reviewer raises a good point. We collated the spatially explicit environmental variables from various sources which we reference. We noticed that the references hadn’t been properly linked in the text. We apologise for this oversight and now provide the full references at the bottom of the supplementary materials. In table S2, we have reworded the explanations for those layers which were developed using the CARS2009 dataset to clarify how these were developed. In combination with the references provided we believe this now provides any interested readers sufficient information. We have provided units for each layer as further context but we have decided not to have BotTemp and SSTGrad layers sequentially because the variables are currently in alphabetical order (we now note this in the table caption).*

Lines 120-125. “The expert assessment focussed on the congruence between predicted taxa distribution and expert view of taxa distribution,.....”. In the supplementary materials for Demersal fish is stated that data come from: “research trawl database ‘TRAWL’ (Niwa, 2014, 2018)”. More details on this data source could be provided in short here (e.g. VMS/Blue boxes for CPUE exc.).

For example, "all catch records were converted into presence" as stated in supp. mat. More than "presence", "temporal persistence" as sustained availability to trawls could be used. As it stands now, it seems that a single event of catch may be converted into "1", but it could be a random phenomenon.

In fact, for Subtidal invertebrates, the data treatment information on fishery-based sampling is more precise: "Lines 35-40" and issues on "opportunistic sampling" (Line 41) are addressed.

*We think that some of the confusion may again be due to the references not being provided in the supplementary materials 1 (we have now ensured they are). The reference for the research trawl database ‘TRAWL’ (Niwa, 2014, 2018) provides links with further information on the sampling, gear types used, etc.*

*The reviewer is correct in that “all catch records were converted into presence”. We then model the suitable habitat of species using these occurrence data following well documented (routinely used) procedures which are detailed (and referenced) in the supplementary materials. For demersal fish and reef fish the consistent sampling methods means that the data represent presence/absence. For benthic invertebrates we need to account for the varied sampling methods which we detail in the supplementary materials.*

*We now clarify this at L100 – 111 in the supplementary materials 1:*

*“To estimate taxonomic distributions, BRT and RF models require locations of both presences (occurrence records) and absences. For biotic groups demersal fish and reef fish, the consistent methods to collect occurrence data means that where a taxon was not sampled, we assume that taxon was absent. For biotic group subtidal invertebrates, the same assumption was made but split by gear type (to account for differences in sampling efficiency). For biotic group macroalgae, given the differences in sampling methods used to collect occurrence data, we used ‘target-group background data’ (Phillips et al., 2009) as absences (referred to here as relative absence), i.e., a location where a different taxon to that being modelled was recorded (Stephenson et al., 2020). In practice different labelling of ‘absences’ does not affect the modelling approach but illustrates differences in the certainty of the absences (and therefore the outputs). Absence (and relative absences for macroalgae) were generated for each taxon from occurrences within taxonomic groups (i.e., demersal fish absences were generated from demersal fish occurrence records). The location of absences and relative absences was required to be at least 1 km from presence data and the number of absence and relative absence data was set to be equal to the number of presences (following best practice outlined in Aiello-Lammens et al. (2015) and Barbet-Massin et al. (2012)).”*

*We have also clarified what we mean by “unique spatial occurrence” in the main body of the article (123 -127) and now point the reader to the supplementary materials where we provide further details.*

*“Taxa records were aggregated spatially to a 1 km grid resolution representing unique spatial occurrences. The number of unique spatial occurrences varied from 50-13926 for demersal fish, 36-339 for reef fish, 70-10804 for subtidal invertebrates, and 50-422 for macroalgae, with a minimum of 30 unique spatial occurrences required for modelling.”*

Line 129. Are authors mean "Figure 2"?

*Yes, thank you. We have changed this now (L137)*

Table 1. Again 235 or 239 (as per Figure 1)?

*Yes, thank you. We have changed this to 235 (L172)*

Figure 3. In which way, a generic user could look directly for a particular species' distribution without considering the proposed taxonomic/ecological groupings (i.e., without downloading the whole set of data)? A Zenodo interface with potential querying for a particular species may increase the appeal of the data set for more generic users. Would it be possible?

*We agree with the reviewer that this added functionality is useful. For specific data queries (i.e., searching for specific taxa / downloading a subset of data) we propose to use the NZ Department of Conservation's online geoportal (L168 – 171). In addition, we now highlight this point and the link to the geoportal in Figure 3 caption (L178 - 181)*

*"Figure 3. The front end of the Zenodo database for the atlas of seabed biodiversity for Aotearoa New Zealand. Also illustrated are the atlas's eight data folders and metadata file as viewed on Zenodo. Image created under standard terms and conditions © Zenodo. Searching and downloading specific taxa (rather than the whole dataset) can be undertaken from the NZ Department of Conservation's online geoportal available at: <https://doc-marine-data-deptconservation.hub.arcgis.com/>"*

*We believe that between the Zenodo and the NZ Department of Conservation's online geoportal we provide an easy means for those interested in downloading the full dataset or to browse and select which layers / taxa to download.*

Lines 183-185. ".....layers provides an accessible source of data layers to inform further research.". Other possible semantic layers of information as maps to be added to the downloading package (to support decision-making) could be already potentially available:

- cargos traffic routes by blue boxes (that determine noise and littering)
- actual trawlers/long-liners operational grounds on continental margins (that determine ghost nets fishing/pollution as discharged nets and wire entanglement; trawl marks reefs damages)
- projected mining areas (future biodiversity treats)

*We thank the reviewer for this suggestion, but we feel that including this information is outside the scope of our study (which is focussed on species' distributions). These data may be available from other sources, including government websites (e.g., fishing activity, management areas, etc) here: <https://maps.mpi.govt.nz/templates/MPIViewer/?appid=96f54e1918554ebbf17f965f0d961e1>*

## **REVIEWER 2**

### General Comments

The dataset documented in this manuscript is a really impressive achievement, and constitutes an extremely valuable resource for the marine biodiversity research and policy communities in Aotearoa NZ and worldwide. The use of data-driven SDMs corroborated with expert judgement has resulted in a really high quality set of species-level maps, and the methodology behind these is generally sound and clearly described. So overall I am very positive about this manuscript. I have a few comments about its structure and content, and some suggestions about maximising access and uptake.

*We thank the reviewer for their constructive comments and overall endorsement of the work. As for reviewer 1, changes have been made to the text following suggestions by both reviewers and line numbers in the response refer to the resubmitted manuscript. All changes to the text are provided as track changes. Detailed responses are provided below.*

### Specific Comments

I like the infographic (fig 1), and hesitate to make suggestions that would overload it with too much information, but I think it would be possible to include a little more info on the different biotic groups - e.g. in addition to number of taxa, maybe also include number (and maybe time-span) of surveys and sampling points?

*We believe that adding this information in the graphic may make the graphic too cluttered. However, we acknowledge the importance of this information and have added this to the figure caption, which now reads (L104 - 111):*

*“Figure 1 Infographic illustrating the development of the atlas described in this study. Presence and absence (p/a) records for four taxonomic groups were combined with over 20 spatially explicit environmental variables and used to run ensembled SDMs of Boosted Regression Trees and Random Forest models. The number of unique spatial occurrences available for modelling ranged between 50-13926 unique spatial occurrences for demersal fish (collected from 1979 – 2016), 36-339 unique spatial occurrences for reef fish (collected from 1986 – 2014), 70-10804 unique spatial occurrences for subtidal invertebrate (collected between 1896 – 2019) and 50-422 unique spatial occurrences for macroalgae taxa (collected between 1896 - 2018). The models were statistically validated using best practice procedures and evaluated by taxonomic experts to further appraise model accuracy. These assessments were incorporated within the metadata of each layer and uploaded, along with the layers themselves, to the Zenodo data portal. World imagery basemap utilised on inset (ESRI 2022).”*

L120-125: I would find it useful here to provide an illustration (maybe in the supplementary material) of different cases - e.g. where model performance statistics and expert assessment agree / disagree particularly well. This would I think help to make this process a bit more intuitive.

*This is a very useful suggestion. We now point to this illustration in the main body of the text (L133 - 134):*

*“(as an illustration, see examples of where model performance statistics and expert assessment agree and disagree in the supplementary materials 1).”*

*and provide these examples in the supplementary materials 1 (L195 – 204):*

*“In most cases, statistical and expert validations were congruent. For example, the predicted distributions of kahawai (*Arripis trutta* – demersal fish), the erect branching deep-water coral genus *Solenosmilia* sp. (subtidal invertebrate); the wrasse *Pseudolabrus miles* (reef fish); and the New Zealand bull kelp (*Durvillaea antarctica*, macroalgae) were all considered to be “accurate” / “very accurate” by experts and had “excellent” statistical model validation scores. In contrast, there were some taxa for which the expert evaluation scores were much lower (“inaccurate” / “somewhat inaccurate”) than the statistical evaluation score (“good”). For example, the pacific salmon (*Oncorhynchus tshawytscha*, demersal fish), the bryozoan Genus *Figularia* (subtidal invertebrates), the wrasse *Coris sandeyeri* (reef fish), the red algae *Champia novae-zelandiae* (macroalgae). Care must be taken when using predictions where there is discrepancy between statistical and expert evaluation (with expert evaluation assumed to be more accurate).”*

Figure 2: There is a reasonably large body of evidence now showing the the rainbow colour palette has a number of shortcomings in terms of accessibility and perceptual biases (for a short recent overview see Westaway 2022 <https://doi.org/10.5194/gc-5-83-2022>). Although the palette used in these maps is not the ‘classic’ rainbow, it shares some of its characteristics, and I would strongly encourage the authors to consider changing their default colour palette, or at least to provide options.

*We acknowledge the reviewers’ legitimate concerns about colour palettes but in this case we feel that the interval classification means that at least perceptual bias may not be an issue. The visualisations of these layers are provided on the geoportal hosted and managed by the New Zealand Department of Conservation and we are not in a position to request these be changed. We believe the colour palette used is adequate for interested parties to explore patterns of distribution and should a more detailed view of the data be warranted we would encourage that they download the layers (from the linked Zenodo database) and use their preferred colour palette.*

L140: The efforts to make the data available are commendable, and for many users the Zenodo archive plus the interactive GIS will be perfectly adequate. As a biodiversity data scientist, however, the thing that would really open up these datasets would be programmatic API access. I do not think this is essential for this release of the dataset, but it is something that I would encourage the authors to consider as a future development (and maybe to discuss that briefly in the ms). As a committed R user I typically look to ROpenSci (<http://ropensci>) for implementations of this kind of functionality, and indeed it appears that the ‘deposits’ package (<https://docs.ropensci.org/deposits/articles/deposits.html>) does already provide means to access data in Zenodo, so it may not be that significant a task to provide some example code or even to package the atlas up - this could really trigger higher uptake among various research communities.

*We thank the reviewer for this suggestion. We agree that this would be a valuable addition and now refer to this as a future development in the L265 - 267. We will explore incorporating this functionality in future iterations of this work.*

*“Furthermore, future iterations of databases may also seek to include programmatic API access to commonly used statistical software (e.g, R statistical software) to facilitate use from biodiversity scientists”*

Supplementary materials

It's not clear from the description of the data, but are the raw occurrence data also deposited in any biodiversity data aggregators (OBIS would be the obvious choice)? If not, and if this is feasible, they would constitute a valuable addition.

*Some of the data are freely available (e.g., demersal fish) whereas others are not (e.g., subtidal invertebrates from the NIWA invert which contains commercially sensitive data). However, the underlying sample data can be available upon reasonable request (now stated L 6 in supplementary materials 1)*

It would be useful to provide a little more information about absences. My understanding is that for both fish datasets and for the invertebrates, absence data is available (i.e., the assumption is that all species from the relevant group were identified at each sampling event, so the absence of a species indicates that it was looked for but not found). Is that correct? And if so, it's not clear why absences needed to be generated using 'target group background data' (from L100) - is it not possible to infer that the surveys where a species was not recorded is an absence, and to use all of these? I may have misunderstood but it would be useful to provide a little clarification. I understand that there is presence-only data for macroalgae - this is clearly explained.

*The reviewer raises a good point as this was not clear in the text. We now elaborate on the different datasets to clarify this at L100:*

*"To estimate taxonomic distributions, BRT and RF models require locations of both presences (occurrence records) and absences. For biotic groups demersal fish and reef fish, the consistent methods to collect occurrence data means that where a taxon was not sampled, we assume that taxon was absent. For biotic group subtidal invertebrates, the same assumption was made but split by gear type (to account for differences in sampling efficiency). For biotic group macroalgae, given the differences in sampling methods used to collect occurrence data, we used 'target-group background data' (Phillips et al., 2009) as absences (referred to here as relative absence), i.e., a location where a different taxon to that being modelled was recorded (Stephenson et al., 2020). In practice different labelling of 'absences' does not affect the modelling approach but illustrates differences in the certainty of the absences (and therefore the outputs). Absence (and relative absences for macroalgae) were generated for each taxon from occurrences within taxonomic groups (i.e., demersal fish absences were generated from demersal fish occurrence records). The location of absences and relative absences was required to be at least 1 km from presence data and the number of absence and relative absence data was set to be equal to the number of presences (following best practice outlined in Aiello-Lammens et al. (2015) and Barbet-Massin et al. (2012))."*

Do you have any thoughts about how reproducible the species distribution modelling process is, given that there is a reasonable amount of model tuning and analyst judgement involved? There have been some efforts to improve the reproducibility of SDMs (e.g. Golding et al. 2017, <https://doi.org/10.1111/2041-210X.12858>) but I am unsure what the current best practice is considered to be.

*We believe that we provide detailed methods which would allow other researchers to reproduce our work given the same data. Furthermore, the model tuning is mostly automated and therefore there wouldn't be much scope for analyst judgment. The idea with these spatial distributions is that they may be periodically updated as and when new data is available. In that sense, these data may "change" every 5 – 10 years but this will be due to differing datasets used in the modelling. In line with this, the expert judgement may vary should different experts be used in future, but in this work*

*we had access to a large number of number experts who arrived at consensus. If this process were repeated, we do not think their assessments would vary greatly.*

#### Technical Comments

Figure 1: Zonodo should be Zenodo

*Changed – thank you*

L129: This should refer to Figure 2, not Figure 1

*Changed - thank you*