

Dear Timothée,

We are particularly pleased and grateful that both reviewers enjoyed our manuscript. We answered all their comments and added a few suggested clarifications. Namely, we clarified that our approach models the fundamental niche of the species, without any factors restricting it as is the case with real-life species; we added details to the modelling section, to the take-home message section and added the true distribution maps alongside the prediction maps in the last Figure. Furthermore, we added a general figure at the start of the paper summarising the workflow of the analysis, following the suggestion of Reviewer 1, which we agree will significantly help the reader navigate the study. Finally, we spotted a mistake in Figure 4 (Figure 3 in the previous version of the manuscript); we addressed it and updated some minor portions of the text accordingly.

We provide detailed responses (in blue) to the reviewers' specific comments below, and attached a tracked changes version of the manuscript (we apologize for the bad quality of figures in this version: we had to compress it to comply the file size limit).

Thank you again for your interest in our study,

Charlotte Lambert and Auriane Virgili

#### Review by [Alejandra Zarzo Arias](#), 02 Mar 2023 12:37

This manuscript explores the differences when building Species Distribution Models for two virtual species using 5 different spatial resolutions of the environmental conditions, and with two different sampling types.

I enjoyed the manuscript a lot, with so many different models and comparisons it is often hard to organize the results, but they appear clear and well summarized. I think the whole manuscript is well written and the ideas are clear, thus, I only have some minor comments:

First of all, I recommend avoiding green+red combination in maps and figures, as color-blind people often cannot distinguish them (e.g., Figure 3, 9, S1-4).

We are always careful in ensuring the accessibility of our work, so we took care to use all the viridis colour palettes (including the "turbo" palette) to ensure that our maps were accessible to colour-blind people. We also double-checked for the accessibility of the figures with a colour-blind colleague.

I think a figure explaining the different sampling methods (segment-based vs. areal-based) would be nice to fully understand the methods.

We agree, and we have added a workflow chart of the analysis carried out in this study to clarify things and help the reader navigate the paper (Figure 1).

Line 4 Abstract: change "wide array OF environmental conditions"

Line 14: "outputS"

Line 48: "distanT"

We incorporated these three suggestions into the text.

Line 49: not a single citation in the entire paragraph. You can give examples of studies using those methodologies.

We added some references of studies in the paragraph.

Line 52: "process" not "processed"

The suggestion has been incorporated, thank you for having spotted the typo.

Line 171: Paragraph "Single variable approach": I believe a further explanation of parameters and distribution of the residuals would be needed for the reader to fully understand model construction.

We added some details into the paragraph to clarify the model parameterisation (Lines 177-182):

*"The models used the number of individuals sighted per sampling unit as response variable. They assumed a Tweedie distribution of residuals and used thin plate regression splines whose complexity was constrained to three inflection points maximum (i.e. four degrees of freedom maximum). This permits the GAM to adjust the complexity of the curve to the data, while avoiding overfitting. The number of individuals sighted per sampling unit was corrected by the area actually sampled within this unit by including this area as an offset in the model ( $2 \times \text{segment length} \times 0.2\text{m}$ )."*

Line 172: the reference should explain GAM models, I think you mean Wood, 2006.

The Wood et al 2011 reference indeed is the most recent reference for the mgcv package, while Wood 2006 is the reference for the overall GAM methodology. We added the mention of the package we used (which was missing), and also added the suggested reference (Line 176).

Line 226: remove A in "highlighted A large"

Line 454: remove WAS in "As was observed"

Line 503: "at which the model will be built"

These three suggestions were incorporated into the text.

For the take home message paragraph I would include more details on what information could be useful and the maximum change in scale admissible to trust some of the results when applying models to the real life (i.e., environmental variables response, suitability maps), which represent very important information for example for conservation purposes, or to select new potential sampling sites.

We expanded the "take-home message" section to derive more detailed advice from the results we present in the study (Lines 523-536):

*"Our simulation study provided clear evidence of stochasticity in the modelling process, thereby urging modellers to caution in fitting models and interpreting resulting outcomes. We also evidenced that classical statistical performance metrics (explained deviances, RMSE. . . ) are not good correlates of predictive quality (for spatial pattern). Despite raster-based modelling being faster in computation (thanks to the lower amount of data points), **segment-based models seemed to be more robust to changes in predictor resolution, and are to be preferred as this approach also avoids** the potential loss of information that could occur when rasterising at scales larger than that of biological significance. However, this pattern is strongly dependent on how much the species distribution and/or environmental conditions are impacted by changes in scale. We therefore advise checking for the resistance of spatial*

*patterns to changes of resolution, both for response and predictor variables, before any analysis. To test for this, it may be considered to quantify the heterogeneity of both the environment and the species distribution, at different resolutions, with tools rooted in landscape ecology. Furthermore, if one decides to opt for raster-based analysis, we strongly recommend carefully checking that the spatial patterns observed with the segmented data are still clearly identifiable after the rasterisation process. Above and foremost, the final choice of sampling type and resolution must depend on the question at hand. This choice must be informed carefully, so that the scale is, as much as feasible, adequately tuned between the observation process and the environmental predictors.”*

Finally, in the main results Figures (also in Supplementary Figures), consider including the original map/graph for the "real" original non-modified models so it makes it easier to compare the results visually. For example, a black curve of the original response in the variable curves graphs.

We added the “true” distribution maps alongside the prediction maps in Figure 10. These are the number of individuals present in each cell for each simulation, then averaged across the hundred simulations (the maps were built for each resolution separately). The text in the corresponding methods and results sections was updated accordingly.

We cannot however include the “true” relationship to environmental condition in the figures in Supplement: these true relationships (Figure 3) were expressed as suitability (rescaled between 0 and 1), while the GAM curves are expressed as abundances. As they are on different unit/scales, we cannot reliably represent them on a single graph. We reworked a bit the Figure 3 however, to make it a bit clearer.

I would recommend the authors to check publications from Vítězslav Moudrý's research group, I believe they will find them interesting.

We did not know of this research group – which is unfortunate, because indeed we have many research interests in common. We added some very relevant references from this group where applicable (namely, Moudrý et al 2023; Lines 38 and 65). Thank you for mentioning it.

#### **Review by anonymous reviewer, 09 Feb 2023 16:11**

This is an interesting and mostly well-executed work. The authors create "virtual species" by defining what is essentially a fundamental niche. They then use a Poisson point-process to generate presences in the region of choice. Their "virtual species" are defined only by their abiotic niche (i.e., biotic interactions and dispersal are disregarded), which is not too much of a problem, although it would be good to be explicit about this: the authors are modeling niches and applying them to geography, rather than modeling the actual biology of a species (including interactions, and movements). Finally, they sample, realistically over the universe of points created by projecting favorable conditions, and proceed to do ecological niche modeling using GAMs, and AIC for goodness of fit. They compare modeling results for different sampling methods, resolutions etc. The results are clearly presented and explained.

I wish the authors refrain from referring to what are projections of niches in geography as "Species Distribution Models" since in reality they are \*potential\* distribution models, but this is mostly a personal preference, since in the literature the distinction between a \*potential\* and an \*actual\* distribution is seldom made. But this is a "virtual species" paper,

and one can only imagine that their species has unlimited dispersal capabilities and no significant biotic limitations. Maybe the authors can state this explicitly in the introduction. We agree with the difference implied behind the two terminologies. We indeed simulated a “simplified” species, for which there is no alteration of the fundamental niche by dispersal barriers, inter-specific interactions or any other factors restricting the realisation of the fundamental niche into the geographical space. Mentioning this explicitly would indeed help the reader in grasping the lessons to be made from this simulation exercise, and we added some details about it into the Method section (Lines 118-120):

*“We simulated two different species per region, whose distributions were driven solely by their fundamental niches. That is, their distributions were driven by environmental forcings alone, without any effects of biological factors like inter- and intra-specific interactions nor dispersal limitations. These environmental forcings were different for each species: one species distribution was driven by...”*

I think the preprint should be recommended, despite the fact that, as with many similar attempts, the complexity of the analysis is great, and one is always left with the question of whether there are general lessons, or the results are artifactual, and how general they are. Nevertheless, the work deserves publication if only because of the stress in the importance of scale of the variables. I believe the authors are basically illustrating that the "Modifiable Areal Unit" problem, that has bothered geographers for the last hundred years, matters when modeling niche-based distributions (see Jelinski DE, Wu J. 1996, The modifiable areal unit problem and implications for landscape ecology. *Landscape ecology*.11:129-40).

We are aware that is one of the main drawbacks of simulations. It is always challenging to understand how to generalise such results to real-life cases. Our results are closely linked to the Modifiable Areal unit problem, yet they offer a different perspective on this, as one of the main objectives was to test for the effect the *type* of sampling unit has on model performance. We were able to test whether a segment-based approach (which is mostly point-based) was as sensitive to the changes of scale as were areal-based approaches. And indeed, we showed that segment-based approaches were more robust, providing the spatial patterns were resistant to changes in scales. To better express the link of our study with this long-standing question, we added “Modifiable Areal Unit problem” as a keyword.