

# 1 Conformal Prediction quantifies the uncertainty of

## 2 Species Distribution Models

3 Timothée Poisot — Département de Sciences Biologiques, Université de Montréal, Montréal QC,  
4 Canada timothee.poisot@umontreal.ca

5 **Abstract:** Providing accurate estimates of uncertainty is key for the analysis, adoption, and  
6 interpretation of species distribution models. In this manuscript, through the analysis of data  
7 from an emblematic North American cryptid, I illustrate how Conformal Prediction allows fast  
8 and informative uncertainty quantification. I discuss how the conformal predictions can be used  
9 to gain more knowledge about the importance of variables in driving presences and absences, and  
10 how they help assess the importance of climatic novelty when projecting the models under future  
11 climate change scenarios.

## 12 Introduction

13 The ability to predict where species may be found is a cornerstone of biogeography and  
14 macroecology (Elith 2019). Techniques from the field of applied machine learning (ML hereafter)  
15 are now routinely used alongside ecological approaches to train generalizable species distribution  
16 models (SDMs hereafter) (Beery et al. 2021). SDMs generate a binary response (corresponding to  
17 the prediction that the species is likely present/absent under given environmental conditions) or a  
18 quantitative score most often as a probability of presence or habitat suitability, indicating how  
19 strongly we believe that the species may be present at the location.

20 Proper communication of the uncertainty associated to the prediction of a SDM is important,  
21 since we usually seek to apply these models to look both forward and backwards in time  
22 (Franklin 2023). This projection if the model to different times is usually called “transfer” (Zurell  
23 et al. 2012), whereby a model trained under historical (baseline) conditions is applied to past/  
24 future projections of the same predictors. The projection of SDMs can also happen in space  
25 (Petitpierre et al. 2016), to predict where species may invade or be naturalized. Even when  
26 predictions are not projected, spatial knowledge of the uncertainty is valuable information: it can  
27 be used to identify areas where the model predictions are trustworthy. Current checklists on the  
28 reproducibility of SDMs emphasize the consequences of data uncertainty (Feng et al. 2019). Yet,  
29 predictions also have inherent uncertainty, which is usually not adequately communicated. This  
30 can be, for example, because of genuine uncertainty about (or inability to capture through the  
31 model) the actual response of the species to combination of predictors (Parker et al. 2024).

32 A common way to capture information about the variability of SDMs is to rely on non-parametric  
33 bootstrapping (Valavi et al. 2021), wherein models trained on random subsets of the data are  
34 compared to estimate the distribution of the response under incomplete sampling. This approach  
35 captures more than one type of variability (Thuiller et al. 2019), and provide valuable information  
36 about the range of performances that can be expected from a model. Other methods are built into  
37 the predictor itself, as is the case for e.g. BARTs (Carlson 2020), which estimate their own  
38 uncertainty. But either situation comes with drawbacks. Bootstrapping requires to train and  
39 evaluate the model hundreds of times, and on partial datasets, which is computationally  
40 inefficient. Using built-in methods limits one to the classifier for which these methods are  
41 available, which prevents for example the use of a new algorithm with the same estimation of  
42 uncertainty.

43 In this manuscript, I illustrate how the ML technique of conformal prediction (CP) allows to  
44 identify instances (combinations of environmental variables) for which a trained and calibrated  
45 model cannot confidently make predictions (Gammerman et al. 1998). A brief introduction to CP  
46 is provided in this manuscript, but the topic is covered in more depth by Shafer & Vovk (2007) for  
47 the mathematical foundations, by Fontana et al. (2020) for a historical perspective, and by

48 Angelopoulos & Bates (2023) for concrete recommendations. By way of contrast to e.g.  
49 bootstrapping, CP does not necessarily involve retraining the same model many times over, but  
50 instead wraps the model into an additional prediction step, and returns estimates of credibility  
51 based on the distribution of past model predictions compared to ground-truthed data. This is an  
52 important difference, as the variability measured through conformal prediction is inherent to the  
53 model, and is not a measure of variability coming through the distribution of data (Lei &  
54 Wasserman 2013). Conformal prediction provides what is essentially (for classification problems)  
55 a confidence interval around the presence or absence of a species in a given location. This is a  
56 particularly important feature, in that CP achieves this in a way that creates several analogues  
57 between ML prediction and fundamental concepts in frequentist statistics (Neyman 1937).

58 One of the reasons why CP is particularly promising for uncertainty quantification in SDMs is  
59 that it is a distribution-free method: it requires neither assumptions about the model nor prior  
60 knowledge of the outcome distribution to provide confidence intervals that are as small as  
61 possible while being *guaranteed* to contain the true value under a set risk level (Vovk et al. 2018).  
62 This is particularly important when transferring a SDM to novel environments (Zurell et al. 2012),  
63 where we expect covariate shift (the joint distributions of predictors are different when training  
64 and predicting), a prediction context that CP is robust to (Fannjiang et al. 2022, Tibshirani et al.  
65 2019).

66 Using occurrence data about an emblematic North American cryptid, I provide a template for the  
67 adoption of CP as a natural way to quantify uncertainty of species distribution models. In  
68 particular, I show how predictions under CP (i) identify areas where the species range is  
69 uncertain, (ii) estimate uncertainty differently from bootstrapping methods, (iii) can be explained  
70 using Shapley values analysis, and (iv) quantify the accumulated uncertainty when transferring  
71 the SDM to future conditions. I conclude by highlighting ways in which using CP can both  
72 simplify the process of training SDMs, and provide information that make their discussion and  
73 analysis more informative.

74    **Methods**

75    *Data*

76    **Occurrence data**

77    The occurrence data used in this article are geo-referenced observations of the Sasquatch (Lozier  
78    et al. 2009). Although these observations are likely to be mis-categorized American black bears  
79    (Foxon 2024), they nevertheless share many features of the data that are used to train SDMs: high  
80    auto-correlation, uneven sampling effort, and clear association with several bioclimatic variables  
81    that is robust enough to train a predictive model. The recorded locations, as well as background  
82    points, are presented in Figure 1.

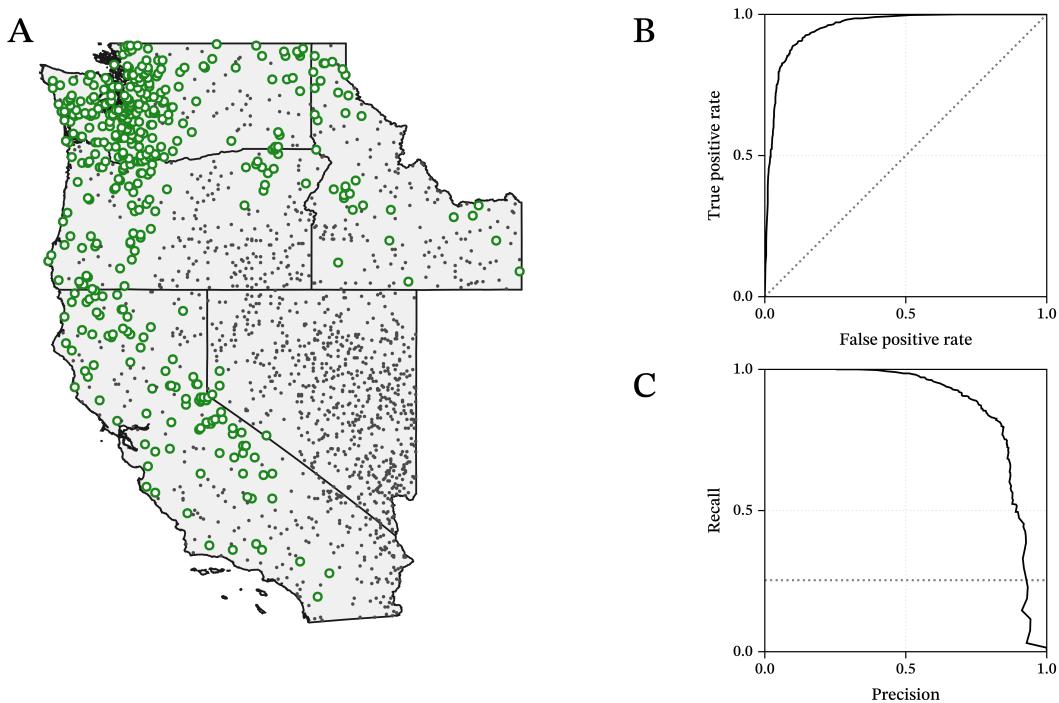
83    **Pseudo-absences generation**

84    The dataset of observations is composed only of presences. In order to establish a baseline of  
85    absences to train a binary classifier, there is a need to generate a number of pseudo-absences,  
86    which simulates locations at which the species, if not absent, has not been observed. In order to  
87    do so, the presence data were first spatially thinned to be limited to one for each cell, at a 5.0  
88    minutes of arc resolution. Cells that had no observation were potential candidates for a pseudo-  
89    absence, and were further selected by drawing a number of them, without replacement, where  
90    the probability of inclusion in the sample was proportional to  $h_{\min}^{-1}$ , where  $h_{\min}$  is the Haversine  
91    (great arc) distance to the nearest cell with an observation, measured in kilometers. In other  
92    words, cells that were close to an observation were unlikely to be included, and cells that were  
93    further away were more likely to be so. To avoid sampling pseudo-absences too close to presences,  
94    the pixels less than 10 kilometers away from known observations were excluded from the  
95    background data.

96 The number of pseudo-absences was arbitrarily set to two times the number of presences.  
97 Although Barbet-Massin et al. (2012) recommend to use the same number of presences and  
98 pseudo-absences for classifiers, using an imbalanced dataset is not a problem: stratified k-folds  
99 cross-validation is perfectly able to handle the moderate class imbalance we introduce  
100 (Szeghalmy & Fazekas 2023), and the model performance (as will be established in a later section)  
101 is sufficient. Moreover, most real-world applications of classification will have to deal with  
102 problems with class imbalance (this is particularly likely to be true of SDM application from  
103 sampling data, where presences may be the minority of outcomes); it is therefore important to  
104 ensure that we do not establish a testing scenario that is too optimistic about the prevalence of  
105 presences. In all cases, class imbalances is a feature of data that must be dealt with in order to get  
106 the more predictive models (Benkendorf et al. 2023).

## 107 **Bioclimatic data**

108 The model was trained, validated, and applied on the 19 WorldClim2 BIOCLIM variables (Fick &  
109 Hijmans 2017), at a spatial resolution of 2.5 minutes of arc. Preliminary analyses using 0.5, 2.5, 5,



110 Figure 1: Overview of the occurrence data (green circles) and the pseudo-absences (grey points) for the  
111 states of, clockwise from the bottom, California, Oregon, Washington, Idaho, and Nevada (A). The  
112 underlying predictor data are at a resolution of 2.5 minutes of arc, and represented in the World Geodetic  
113 System 1984 CRS (EPSG 4326). The panels on the right column show the ROC curve (B) and PR curve (C),  
114 with the random classifier indicated by a dotted line. The area under the ROC curve is  $\approx 96\%$ .

115 and 10 minutes of arc show that the qualitative results presented hold. For the projection of the  
116 model under climate change, I only report the future data under the SSP370 scenario (“business  
117 as usual”), for the MRI ESM2-0 GCM, over the period 2081-2100.

118 The climatic novelty of the baseline *v.* future data is estimated through the Euclidean distance  
119 (Fitzpatrick et al. 2018), specifically by assigning as a novelty score for each pixel in the future the  
120 distance to its closest baseline analogue. This novelty is measured on de-meanned predictors with  
121 unit variance.

## 122 *Species distribution model*

123 All analyses are conducted using the `SpeciesDistributionToolkit` package (Poisot et al. 2025) for  
124 *Julia* 1.11.

## 125 **Model structure**

126 The model used here is a logistic regression, with interactions terms up to a maximum degree of  
127 two (preliminary analyses with random forests, naive Bayes classifiers, and rotation forests gave  
128 similar results). When trained on a vector of features  $\mathbf{x}_i$  (with null means and unit variances), the  
129 model will return a probability  $p_+$ , which correspond to the probability of these environmental  
130 conditions being associated to the presenceof the species. This probability is turned into a  
131 presence/absence decision by comparing it to a threshold, as explained in a later section. Because  
132 this logistic regression is a deterministic classifier, the prediction  $p_i +$  statisfies  $0 \leq p_i + \leq 1$ , and  
133 we use  $p_- = 1 - p_+$  as the probability that the species is absent from the location.

## 134 **Tuning**

135 We tune this model by (i) iteratively forward selecting the best set of predictor variables, and (ii)  
136 optimizing the threshold  $\tau$  above which a site with a probability for the positive class  $p_+$  is  
137 considered to be positive (turning the prediction of presence into  $p_+ \geq \tau$ ). In both cases, the  
138 cross-validation strategy is the same: the dataset is split in 10 random folds, 9 of which are used  
139 for training and one for validation. All folds are used for evaluation, providing exhaustive cross-  
140 validation. The folds are stratified so that the relative number of present cases in the training set

141 is similar to that of the entire dataset. The performance on each set, for the purpose of defining  
142 the set of variables to include of the threshold to use, is measured as the average of the Matthews  
143 Correlation Coefficient (MCC) across each of the ten folds. The MCC is the most accurate  
144 representation of a binary classifier performance (Chicco & Jurman 2023), and avoids the pitfalls  
145 of several other validation measures.

146 For all steps of model training and validation, the identity of instances composing the different  
147 folds remains fixed. This ensure that the changes in MCC are only due to the addition of the  
148 variable, and not to the random sampling of a training/validation set with different properties.  
149 Although some authors encourage the use of spatially-stratified cross-validation (Soley-Guardia  
150 et al. 2024), this is not a desirable strategy for this use-case. The area in which the predictions will  
151 be made is entirely delimited by the bounding box of observed presences, and there is therefore  
152 no risk of covariate shift when shifting from validation to prediction (outside of the situation of  
153 temporal transfer of the SDM).

154 The predictors included in the model have been decided through the use of forward selection.  
155 This is an important step in order to perform dimensionality reduction (which generally increases  
156 the predictive accuracy), but also to ensure that the set of retained variables is reduced enough  
157 that it can be interpreted. Variables were retained as part of the final set of predictors if adding  
158 them increased the MCC for the model once retrained with this new variable.

159 One of the most efficient ways to increase the performance of binary classifiers is to change the  
160 decision rule leading to a positive (here, presence) prediction, so that presences are assigned  
161 when  $p_+ \geq \tau$  – a process known as moving threshold classification (Liu et al. 2013, 2016). The  
162 value of  $\tau$  is an hyper-parameter of the model, which is chosen to maximize the value of a  
163 measure of model performance (here the MCC) when evaluated over many different values. In  
164 this instance, we optimized the value of  $\tau$  by picking the value out of 200 linearly spaced value  
165 between the smallest and largest prediction made on the training set. The value of  $\tau$  that  
166 maximizes the MCC during cross-validation was selected as the optimal threshold for the  
167 classifier. Note that even though our decision rule for the presence of the species is  $p_+ \geq \tau$ , we  
168 will keep the information about  $p_-$  as is required for conformal prediction.

169 **Bootstrap variability**

170 Bagging (bootstrap aggregating) is often used as a measure of uncertainty to the underlying data  
171 when training SDMs (Beale & Lennon 2012). When performing bagging, the model is trained on  
172 samples drawn with replacement from the training set (which leaves out approx. 37% of the  
173 dataset). Models are then evaluated on samples that were not used as part of their training,  
174 usually using cross-validation (Bylander 2002) or measures of the out-of-bag error (Janitz &  
175 Hornung 2018). Although ensemble models *can* result in a better predictive performance  
176 compared to single models (Drake 2014), this is not a guarantee (and depends on the structure of  
177 the bias/variance trade-off for the specific model and its training set). The many models trained  
178 on the bagging dataset form an homogeneous ensemble, which is to say a set of models that share  
179 the same algorithm and hyper-parameters, and only make different predictions as the result of  
180 having been trained on different subsets of the full training set.

181 Measures of whether the different models composing the homogeneous ensemble agree can  
182 provide a measure of the effect of data and parameter uncertainty (Petropoulos et al. 2018), or  
183 what Davies et al. (2023) termed the “SDM uncertainty”. The best model identified after  
184 thresholding was evaluated on a hundred bootstrap samples, yielding an homogeneous ensemble  
185 model from which we estimate bootstrap variability (Chen et al. 2019). Because the model is kept  
186 constant in this analysis, the measure of variability we will derive from the ensemble model is an  
187 estimate of how sensitive the estimation of the model parameters is to small perturbations  
188 (specifically: spatially homogeneous under-sampling) to the training data.

189 **An introduction to conformal prediction**

190 Conformal prediction differs from regular prediction in that, rather than a single point prediction,  
191 it returns sets corresponding to the ensemble of *credible* outcomes given an input  $x$  representing  
192 environmental conditions at which we seek to make the prediction. Given the observed quantiles  
193 of the model output on validation data, these sets are obtained through a simple calibration step.

194 Therefore, CP requires an already trained model, and is agnostic to the process through which  
195 this model is trained. In this section, I highlight two important features of CP: the notion of  
196 *credible sets* (and how they are obtained), and the notion of *coverage*, which is a measure of  
197 tolerance to error.

## 198 *Understanding conformal predictions*

199 By contrast to the non-conformal SDM, the conformal classifier returns, for an input of  
200 environmental predictors  $\mathbf{x}$ , a set  $C$  containing the “credible outcomes” for this prediction. This  
201 set is termed the *credible set*, and under a binary classification task (the species is either present or  
202 absent), there are four possible combinations for the content of credible sets:  $C = \{+\}$ ,  $C = \{-\}$ ,  
203  $C = \{+,-\}$ , and  $C = \emptyset$ .

204 The first two cases are simple: if the credible set contains a single output, the model can  
205 confidently make a prediction that excludes the other class. In the case of  $C = \{+\}$ , for example,  
206 the point prediction for the presence score  $p_+$  is high enough that the outcome of absence can be  
207 ruled out given the known predictions on training examples. In some cases, the credible set may  
208 contain both classes, as in  $C = \{+,-\}$ . Although they may not be *equally likely* (there is no  
209 guarantee that  $p_+ \approx p_-$ ), the scores are close enough to not confidently exclude one of the  
210 outcomes from the model prediction. In the specific cases of SDMs, these correspond to areas of  
211 true uncertainty, where the known training examples credibly support both the presence or  
212 absence of the species. The final situation,  $C = \emptyset$ , corresponds to pathological cases where  
213 *neither* outcome can be credibly supported. Given the training data (and the distribution of  
214 presences and absences), the model is not able to make a prediction for this input. The increased  
215 frequency of such predictions is most likely a strong sign that the risk level is too high (the  
216 confidence interval is too broad) for the training data given to the conformal model.

217 These situations correspond to four different outcomes in terms of the SDM certainty about the  
218 distribution of the species. The most intuitive situation is  $C = \{+\}$  or  $C = \{-\}$ , in which case the  
219 conformal model predicts that the absence (resp. presence) of the species is *not* a credible  
220 outcome for the environmental conditions given as an input. Throughout this manuscript, I will

221 refer to these predictions as “sure presences” and “sure absences”, as they convey the information  
222 that there is no reason to expect that the prediction is uncertain. The second situation,  $C = \{+, -\}$ ,  
223 corresponds to inputs for which the presence and the absence of the species are credible, and I  
224 will refer to them as “unsure”. The rare cases where  $C = \emptyset$  will be “undetermined” predictions.

225 *Obtaining conformal predictions*

226 There are several ways to decide whether a point prediction from the model results in which  
227 credible set. A core assumption of CP is that the data used for training should be exchangeable, or  
228 in other words, their joint probability distribution should be (close to) invariant under finite  
229 permutations (Aldous 1985). This will almost never be the case for data with a spatial structure;  
230 nevertheless, this does not rule out the use of CP for species distribution modeling, as Oliveira et  
231 al. (2024) show that CP is acceptably robust to lack of exchangeability.  
232 The central idea of CP is to associate a conformal score to a point prediction. This can be achieved  
233 by applying the softmax function to the values for  $p_+$  and  $p_-$ , giving

$$s_+ = \frac{\exp p_+}{\exp p_+ + \exp(1 - p_+)}, s_- = \frac{\exp(1 - p_+)}{\exp p_+ + \exp(1 - p_+)} \quad (1)$$

234 The conformal score associated to a prediction is  $1 - s_\cdot$ , where  $\cdot$  is the prediction (+ or -) made  
235 by the model. We call the distribution of conformal scores  $\mathcal{S}$ . Note that this can be done without  
236 using the softmax function, but it is included here as it is best practice for classification.  
237 The next step is to identify a critical value  $\hat{q}$  above which a conformal score indicates that the  
238 prediction it describes is credible. This critical value is picked by examining the empirical  
239 quantile distribution of the conformal scores calculated over  $n$  training examples, and an  
240 acceptable level of risk  $\alpha$  (explained in depth in the next sub-section), and specifically by  
241 identifying the  $q_i$ -th quantile, where

$$q_i = \frac{[(n + 1)(1 - \alpha)]}{n} \quad (2)$$

242 The corresponding value of  $S$  below which a proportion  $q_i$  of values lies is  $\hat{q}$ . In other, more  
243 intuitive words, the value  $q_i$  indicates what proportion of wrong classification events we must  
244 accept before we have accumulated enough evidence to be confident about a prediction. When  
245 performing the prediction, we calculate the score of a new prediction according to Equation 1.  
246 For every possible class  $x$ , if  $s_x \geq (1 - \hat{q})$ , this class is retained as part of the credible set.  
247 The value of  $\hat{q}$  can be obtained either through using a holdout set for training (Split Conformal  
248 Prediction), by retraining the model in a way akin to Leave-One-Out cross-validation (Full  
249 Conformal Prediction), through the use of quantile regression (Romano et al. 2019), or through  
250 taking the median of several estimates of  $\hat{q}$  after cross-validation (Vovk et al. 2018). In this  
251 manuscript, I employ the later method, as it provides a rapid and statistically acceptable estimate  
252 of  $\hat{q}$ , without requiring too much computing time.  
253 To summarize, the output of the conformal classifier is, in a sense, a point estimate of the credible  
254 outcomes of a model, using the value estimated for  $p_+$  as well as knowledge about which of these  
255 were associated to the correct label in the training data. A location is defined as included in the  
256 range if the positive outcome is included within the credible set returned by the conformal  
257 classifier, and as excluded from the range when it is not. Because the conformal classifier can  
258 identify that both outcomes are credible based on the training data (while giving them different  
259 weights), predictions in which both the positive and negative outcomes are included in the  
260 credible set can be seen as “uncertain” at this given risk level.  
261 How frequently a specific prediction is uncertain is termed the inefficiency of the classifier, which  
262 is defined as the average cardinality of all credible sets. The inefficiency is bounded upwards by  
263 the number of classes (two for binary classification); when the inefficiency is  $\approx 1$ , the conformal  
264 classifier behaves (essentially) like deterministic classifier, by returning a single class for each  
265 instance. An inefficiency close to unity is not desirable: smaller sets can hide our actual  
266 uncertainty (Sadinle et al. 2018). Because the conformal models wraps the logistic regression  
267 model, we can further divide the “unsure” predictions as a function of whether they would be  
268 within the range as predicted by the SDM, which I will call “unsure presences”; the other unsure  
269 predictions are referred to as “unsure absences”.

270 *The coverage level*

271 CP allows users to set a desired error rate,  $\alpha$ , which appeared in Equation 2. Intuitively, what CP  
272 does, is inform the user on whether the credible set contains the true value with probability  $1 - \alpha$ ,  
273 which allows to directly interpret this value as a true confidence interval. This error rate is usually  
274 referred to as the *marginal coverage*, in that it captures the probability of success marginalized  
275 over the known validation points. Because the estimate of uncertainty involves the original  
276 model, it is important to apply CP on a model with adequate performance.

277 Chaning the risk level  $\alpha$  leads to different estimates of how commonly multiple classes will be  
278 accepted as a credible outcome. Using a low level of risk ( $\alpha \approx 0$ ) yields usually leads to all  
279 outcomes being credible ( $\hat{q} \approx 1$ ), at the cost of a very high uncertainty. When values of  $\alpha$  get too  
280 large ( $\hat{q} \approx 0$ ), no class can be confidently predicted, and the model will eventually always return  
281  $C = \emptyset$ . Although this later situation is more difficult to make sense of intuitively, a value of  
282 inefficiency that gets smaller than unity should be interpreted as a model that accumulates more  
283 uncertainty (at a given risk level) than the data can support (Romano et al. 2020). Conformal  
284 prediction can therefore inform us on the acceptable risk levels we can operate under given a  
285 trained predictive model.

286 In the rest of this analysis, I will set  $\alpha = 0.05$ . As noted by Angelopoulos & Bates (2023), this  
287 corresponds to estimating whether a specific prediction falls within, or outside of, the 95%  
288 confidence interval across all predictions, which is a convenient callback to frequentist statistics'  
289 usual risk tolerance. Recall that the CP credible sets are estimated based on the model output, and  
290 therefore even when aiming for full coverage, there may be non-ambiguous combinations of  
291 environmental predictors.

<b>Measure</b>	<b>Validation</b>	<b>Training</b>	<b>Ensemble</b>
MCC	0.75	0.76	0.76
NPV	0.93	0.93	0.94
PPV	0.82	0.83	0.82
$\kappa$	0.75	0.76	0.76
TSS	0.74	0.75	0.76
Accuracy	0.91	0.91	0.91

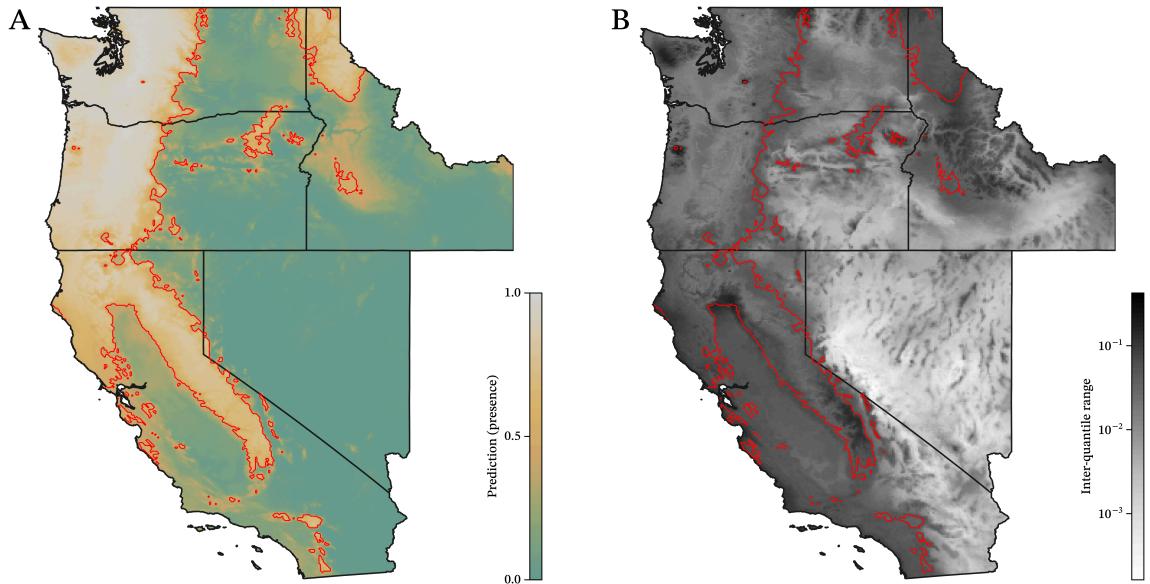
292  
293  
294  
295  
296  
297  
298  
299 Table 1: Overview of measures of model performance for the validation and training sets of the SDM, as  
300 well as the same measures for the ensemble model (measured on the out-of-bag models only). The values of  
301  $\kappa$  and the true-skill statistic are generally comparable to the MCC, but are included as they are commonly  
302 reported in the SDM litterature (Allouche et al. 2006). The high values of the negative and positive  
303 predictive values indicate that the model is suitable to detect both presences and absences.

## 304 Results

### 305 *Performance of the baseline model*

306 In panels B and C of Figure 1, we report the ROC and PR curves for the model. As evidenced by  
307 both these diagnostic tools, the model achieves a very high predictive accuracy. In Table 1, we  
308 report additional measures of performance for the training and validation set of the model (so as  
309 to ensure that the model is not performing better on training data), as well as a measure of the  
310 performance of the ensemble, to show that it can make valid predictions in addition to  
311 quantifying variability. These results confirm that the model is able to identify areas that are  
312 suitable to the species, and can be used for CP.

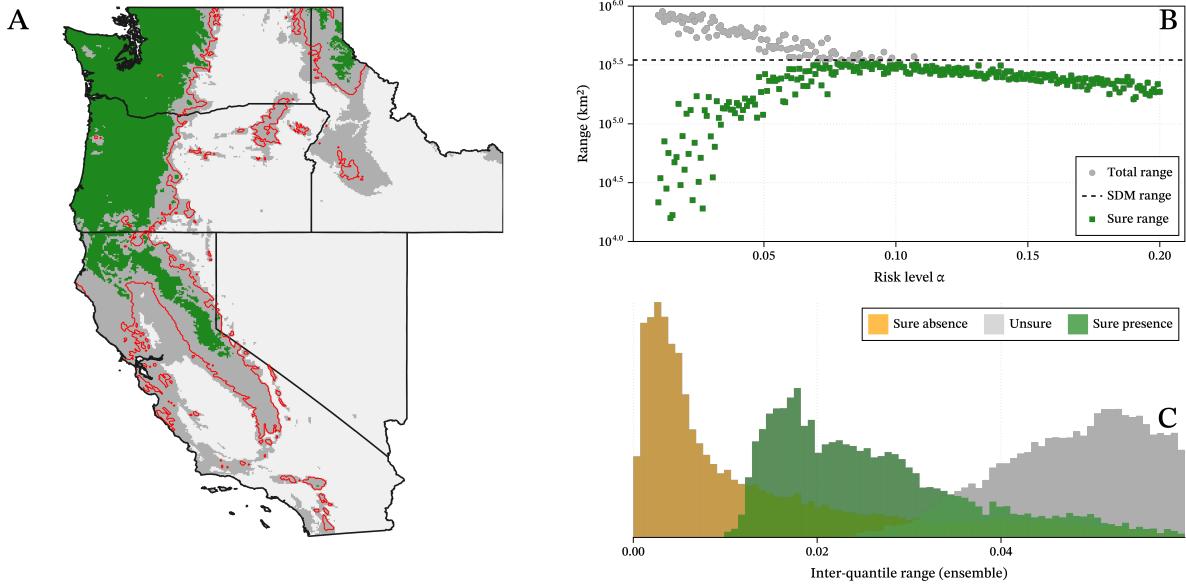
313 Before applying CP, it is useful to examine the output of the SDM in space. The predictions of the  
314 model for the entire region are given in Figure 2, alongside information about the model  
315 variability. Areas of lowest variability (according to the IQR based on non-parametric bootstrap  
316 results from the ensemble) seem to be associated with the absence of the species, with the  
317 variability mostly increasing within the predicted range.



318 Figure 2: Overview of the probability  $p_+$  returned by the model (A), and the inter-quantile range of the  
 319 non-parametric bootstrap model predictions (B). The range, i.e. the limit of cells for which  $p_+ \geq \tau$ , is  
 320 indicated by a solid red line; I maintain this convention for all subsequent figures. Note that the scale of the  
 321 variability is logarithmic, as the model shows good performance and therefore has low variability overall.

322 *Conformal prediction of the species range*

323 Before discussing the spatial output of running the conformal model, it is worth considering why  
 324 the thresholding step as visualized in Figure 2 is not really providing us with a set of certain  
 325 presences and absences. When optimizing the threshold  $\tau$  above which a prediction  $p_+$  from the  
 326 non-conformal model is determined to be a presence, we inherently establish a sort of certain  
 327 presences and certain absences, specifically by ignoring the possibility that there can be uncertain  
 328 predictions. Indeed, the space covered by positive predictions is usually interpreted as the  
 329 (potential) distribution of the species. But this prediction conveys a false sense of certainty, that  
 330 has to do with the very nature of the threshold we optimize. By definition, the threshold is the  
 331 value that finds the best balance between the false/true positive/negative cases on the validation  
 332 data; this is in fact why the optimal threshold is the point closest to the corners of the ROC and  
 333 PR curves indicating a perfect classifier (Balayla 2020). When a prediction  $p_+$  gets closer to the  
 334 threshold, a small perturbation to the environmental conditions locally could bring it on the other  
 335 side of the threshold, and therefore flip the predicted class using the non-conformal classifier.  
 336 Around the threshold is where we expect uncertainty to be the greatest.



337 Figure 3: Overview of areas where the presence of the species is certain according to the CP model under a  
 338 risk level  $\alpha = 0.05$  (A). The certain areas are in dark green, and the uncertain areas, wherein both presence  
 339 and absence are credible, are in dark grey. (B) Surface covered by the sure absence and total range  
 340 (including the superficy of the unsure area) for different risk levels. Note that for  $\alpha \approx 0.1$ , the total  
 341 predicted range starts being lower than the range predicted by the SDM, and the uncertain range collapses.  
 342 (C) Distribution of variability from Figure 2B by type of CP model outcome.

343 To bring these considerations into a spatial context: we expect the areas where the score for the  
 344 present class are closer to the threshold (the limits of the predicted range of the species) to be the  
 345 most uncertain. Importantly, this is true *both* for areas that are inside the range (for which  $p_+$  is  
 346 just above the threshold) and for areas that are outside of it (for which  $p_+$  is just below the  
 347 threshold). CP is perfectly suited to solving this issue, by identifying the areas where one class is  
 348 predicted, but the other class is also credible. In this section, we will project the areas with  
 349 uncertain predictions, and compare the uncertainty quantified by the conformal model to the  
 350 uncertainty derived from the ensemble model.

351 In Figure 3, we show that this prediction indeed stands: the range as predicted by the SDM  
 352 (fig. 3A) falls within the range of unsure predictions. We also see that lowering the risk level  $\alpha$   
 353 leads to a contraction of the area (in km<sup>2</sup>) considered to be credibly associated to only the  
 354 presence of the species ( $C = \{+\}$ ), while the range that is ambiguous ( $C = \{+, -\}$ ) increases  
 355 (Figure 3B). As far as ecologists are concerned, the areas in which the credible set only has a score  
 356 for the absence of the species are the easiest to make sense of: they correspond to regions where  
 357 the model is certain (under the specified risk level) that the species is absent. All other areas

358 (assuming that there are no predictions for which the credible set is empty, which I discuss in the  
359 next section) are *potentially* part of the range of the species: some certainly, some uncertainly.  
360 Depending on the purpose for which the SDM is produced, the uncertain areas can be treated  
361 differently. As Prescott et al. (2025) argue, when dealing with invasive species, it may be more  
362 reasonable to interpret SDMs by erring on the side of caution, which here would mean  
363 considering that unsure presence area should be considered part of the species's range. On the  
364 other hand, when SDMs are meant to guide conservation actions that are costly or should be  
365 focused on areas of high certainty of suitability for the target species (Pěknicová & Berchová-  
366 Bímová 2016), it may make sense to ignore the unsure presences.

### 367 *Relationship between variability and uncertainty*

368 Note that the relationship between the certainty associated to CP, and the variability under the  
369 ensemble model presented in Figure 2B is nuanced: in fig. 3C, it appears that although areas  
370 identified as unsure using CP tend to have higher variability, there is considerable overlap  
371 between the categories. Intriguingly, the overlap between areas that are uncertain according to  
372 the conformal classifier, and areas that are uncertain according to the bootstrap model, is  
373 imperfect. There are a number of points classified as sure presences for which the IQR is very  
374 high, **i.e.** points whose certainty is not affected by undersampling the training data. Notably, the  
375 results in fig. 3C show that it is not possible to find a cutoff in the measure of bootstrap variability  
376 that would identify areas of model uncertainty. This suggests that the classification of predictions  
377 as certain/uncertain according to the conformal prediction is in part reflecting genuine  
378 uncertainty in the underlying data, but also contributing novel information about the fact that  
379 some instances are more difficult to call.

380 These results can be better understood by contrasting what “uncertain” means in the context of  
381 CP, and how it differs from the uncertainty in the ensemble model. The uncertainty derived from  
382 the ensemble model represents whether many models trained on small perturbations of the full  
383 training dataset would agree on a specific prediction task, represented by an array of  
384 environmental predictors. Therefore, the uncertainty from the ensemble originates in the

385 estimation of the parameters, and its sensitivity to being able to access the full information within  
386 the training data. Uncertainty in the conformal classifier is coming from comparing the  
387 prediction to all other predictions under an estimation of the distributions for the conditions  
388 leading to the prediction of the presence (or absence) outcome. Therefore, the uncertainty from  
389 the conformal predictors accounts for all the predictions the model can make, and accounts for  
390 the variability *across* predictions within a fully known dataset.

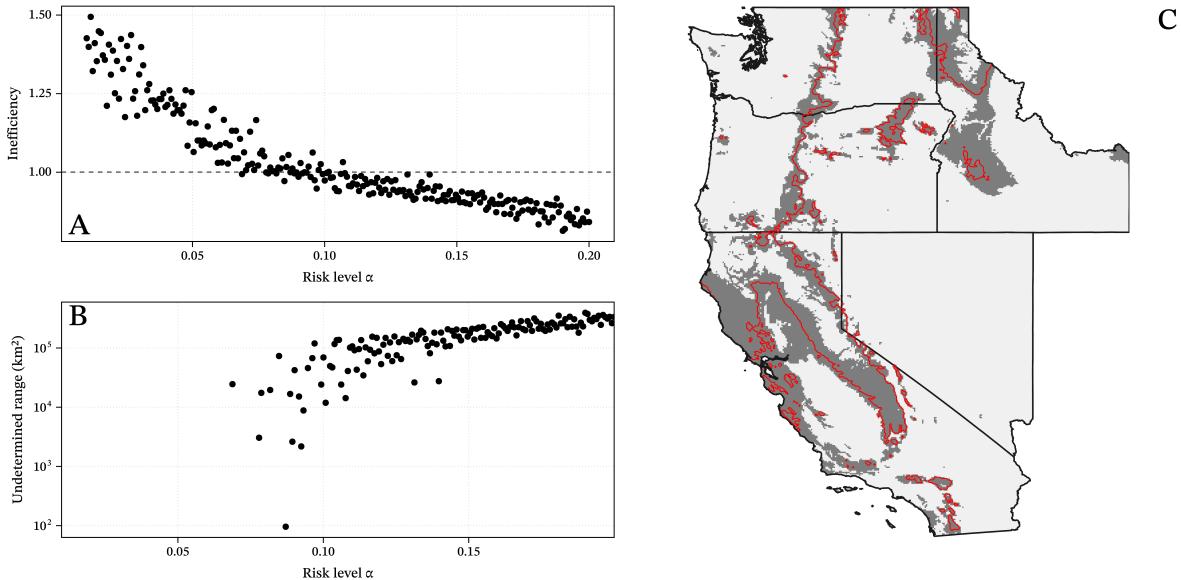
### 391 *Identification of undetermined areas*

392 In Figure 3B, we see that there is a risk level above which the total predicted range starts to get  
393 lower than the range predicted by the SDM. We can explain this behavior through the lens of the  
394 number of undetermined predictions, *i.e.* the number of inputs for which the CP model returns  
395  $C = \emptyset$ .

396 In fig. 4A, we see that above  $\alpha \approx 0.1$ , the inefficiency of the classifier starts to fall under 1 - this  
397 indicates that *on average*, the model is returning fewer than one output for each prediction. In a  
398 sense, this creates an upper limit to the risk we can accept: the model trained on this dataset does  
399 not support conformal prediction for larger risk levels. In fig. 4B, we see that this change of  
400 behavior in the model is indeed resulting in an increase in the range for which the model makes  
401 no prediction, which gets larger when the risk level is too high. The spatial distribution of  
402 undetermined areas is shown in fig. 4C for  $\alpha = 0.2$ : these areas are concentrated around the range  
403 limit as identified by the SDM. This suggests that using a risk level that is too high would result to  
404 no conformal predictions being made for the areas where our need to accurately quantify  
405 uncertainty are the most important, and calls for a cautious investigation of the appropriate risk  
406 level.

### 407 *Model explanation*

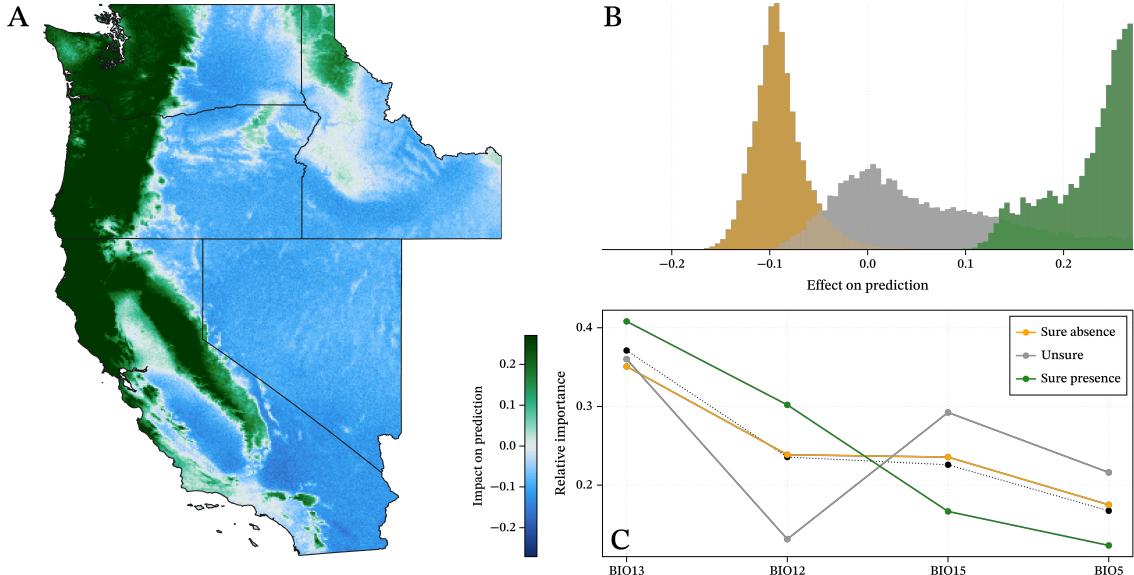
408 In this section, I perform an analysis of Shapley values of the conformal predictor, in order to (i)  
409 assess the importance of variables and (ii) provide explainable results about the relationships  
410 between predictors and response. I rely on the common Monte-Carlo approximation of Shapley



411     Figure 4: Inefficiency (average number of classes in the credible set) for various levels of  $\alpha$  (A); above  $\alpha \approx$   
 412     0.1, the conformal prediction starts returning empty credible sets. This results in an increase in the spatial  
 413     area for which no prediction can be made (B). For  $\alpha = 0.2$ , these areas are distributed around the limit of  
 414     the predicted range, showing that the areas in which uncertainty quantification are most important cannot  
 415     be predicted.

416     values (Roth 1988, Touati et al. 2021). Monte-Carlo Shapley values represent, for each prediction,  
 417     how much the  $i$ th variable contributed to moving the prediction away from the average  
 418     prediction. The Shapley value associated to variable  $i$  is  $\varphi_i \in [-1, 1]$ , which measures how much  
 419     this variable modified the *average* prediction for this class. Shapley values have a number of  
 420     desirable properties regarding the explanation of prediction of responses for environmental  
 421     studies (Wadoux et al. 2023), including their additivity: for any given prediction,  $p = \hat{p} +$   
 422      $\sum_i^{\text{variables}} \varphi_i$ . Because of this additive property, the importance of variables across many  
 423     predictions is usually measured as the average of  $|\varphi|$ , where both positive (the class is more  
 424     likely) and negative (the class is less likely) are counted. This measure of variable importance  
 425     represents the relative impact that each variable had on the process of moving all predictions  
 426     away from the average prediction and towards its actual value. Because Shapley values are both  
 427     additive and independent, they can be measured and aggregated for any arbitrary stratification of  
 428     the data (which allows reporting them conditional on the uncertainty status of the prediction).

429     As the predictions of the conformal model can be split by whether they are certain or uncertain,  
 430     they offer a unique opportunity to delve into the mechanisms that *generate* this uncertainty.  
 431     Namely, if the relative importance of variables is different across these classes of predictions, this



432 Figure 5: Overview of the effect of the most important predictor (A); areas with high values indicate that  
 433 the value of BIO13 at this location make the presence of the species more likely. These values are associated  
 434 to different prediction certainties (B), with predictions within the unsure range being centered around 0  
 435 (i.e. not moving the needle on the average prediction one way or another). Nevertheless, the contribution of  
 436 the variables in different uncertainty categories are different (C), suggesting that Shapley values can help  
 437 create explanations of where uncertainty originates.

438 is strongly suggestive of the fact that there are certain environmental conditions (represented by  
 439 combination of values for each variables) that create or reduce uncertainty. Furthermore, because  
 440 we can split the certain predictions into a presence and absence class, this is a unique opportunity  
 441 to investigate whether the factors leading to a species being present or absent are the same. An  
 442 example of the spatial contribution of a variable is given in Figure 5A.

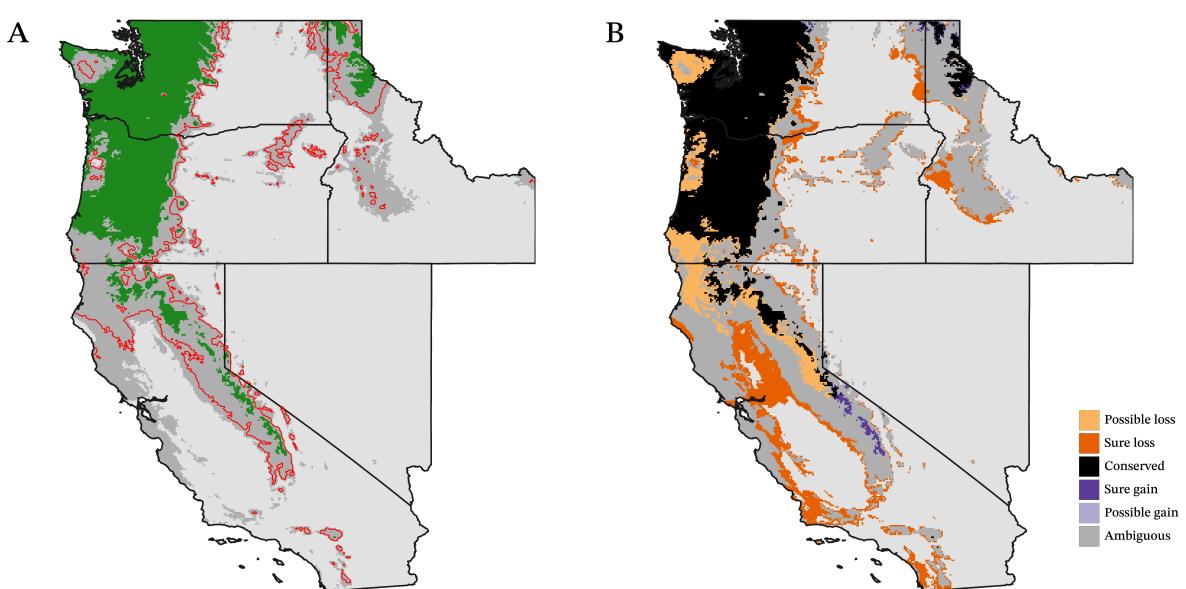
443 We find that, for the most important variable (*i.e.* the one with the largest  $\sum|\varphi|$ ), the contribution  
 444 of this variable tracks the status of the prediction: it tends to be negative when the absence is  
 445 certain, positive when the presence is certain, and around zero when the prediction is unsure  
 446 (fig. 5B). This is a fairly remarkable result, in that it ties Shapley values (a tool to help with ML  
 447 models interpretation) to CP (a technique to accurately convey uncertainty). In Figure 5C, I  
 448 present the relative contribution of all selected variables split by the status of the prediction; this  
 449 reveals that the Shapley values for sure presences and unsure areas are distributed in different  
 450 ways. Notably, BIO15 is far more important in areas of high model uncertainty than in areas of  
 451 either sure presences or absences. This suggests that the division of the prediction according to  
 452 CP status can provide information about which sets of environmental conditions are driving the

453 uncertainty, thereby providing useful information to guide future sampling or model  
454 interpretation.

455 ***Conformal prediction and climate-induced range shifts***

456 In a recent contribution, Smith & Levine (2025) suggest that because of issues around the use of  
457 thresholds, projections of SDMs under climate change scenarios may benefit from a more  
458 continuous perspective. In this section, I present a comparison of the conformal prediction of the  
459 range under a climate change scenario (SSP370. 2081-2100), to illustrate how the future  
460 conformal range can convey information about the certainty of some types of range shift. These  
461 results are presented in Figure 6.

462 Based on the comparison between the baseline (fig. 2A) and projected (fig. 6A) ranges, we can  
463 establish a series of transitions and their interpretations as range change scenarios, which are  
464 presented in fig. 6B. Areas that are certain both now and in the future,  $\{+\} \rightarrow \{+\}$ , can safely be  
465 assumed to be conserved. Areas that where unsure and become surely negative,  $\{+,-\} \rightarrow \{-\}$  are  
466 *possible losses*, as they may have been presences in the baseline data, but are considered lost in  
467 the future. The reverse scenario,  $\{+,-\} \rightarrow \{+\}$ , corresponds to *possible gains*. Sure losses of range



468 Figure 6: Overview of the conformal prediction of the range for the future climate data, equivalent to  
469 fig. 2A (panel A). Spatial distribution of areas where loss and gain are expected to be possible v. certain, as  
470 explained in main text (B).

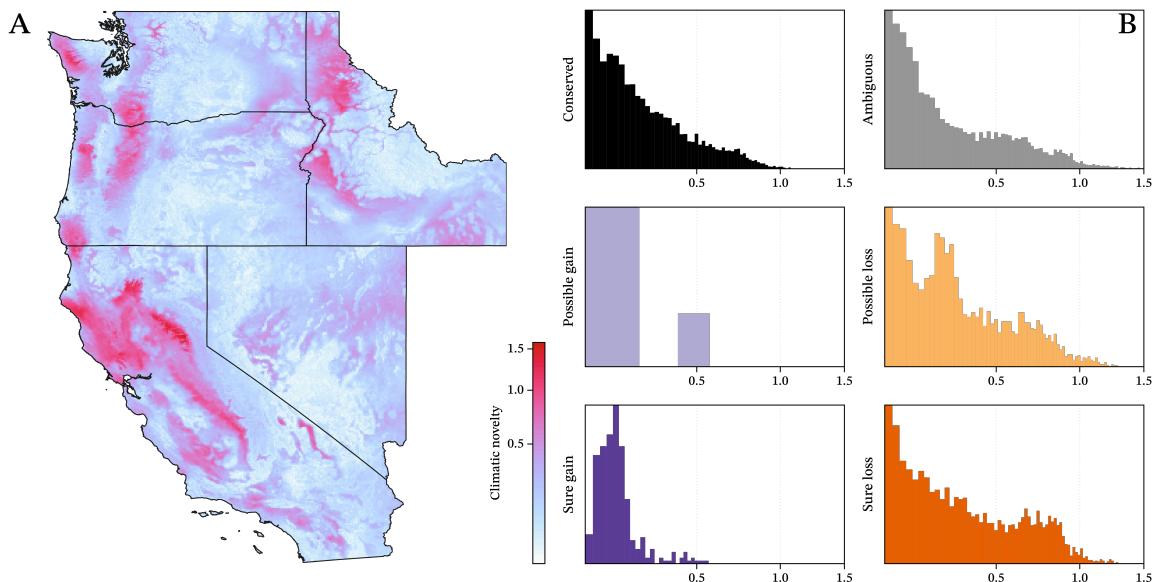
471 correspond to the transition  $\{+\} \rightarrow \{-\}$ , and sure gains of range correspond to  $\{-\} \rightarrow \{+\}$ . Other  
472 situations are considered ambiguous.

473 By applying these rules on the predicted changes in presence/absence status, we can identify  
474 large areas that are confidently loss towards the Southern edge of the species's range, with very  
475 limited areas of either possible or sure gain, strongly suggesting that this species would undergo  
476 range contraction. Note that the area corresponding to ambiguous transitions is relatively large,  
477 which provides a good understanding of the possible variation to be expected under this climate  
478 change scenario.

479 *Conformal prediction and climatic novelty*

480 Zurell et al. (2012) highlight the importance of fully considering uncertainty when transferring  
481 the model to novel climate data: there is a chance that the future climate conditions will not have  
482 occurred in the training dataset, and therefore our confidence in the model outcome should be  
483 lowered. This covariate shift is well documented to decrease the performance of models  
484 (Mesgaran et al. 2014), and CP offers an opportunity to shine a different light on this  
485 phenomenon.

486 This task is particularly crucial given that entirely novel climatic conditions are likely to become  
487 the norm (Mahony et al. 2017), which in turn will drive the emergence of a novel biosphere  
488 globally (Kerr et al. 2025, Ordonez et al. 2024). In this section, I compare the results of  
489 conformation prediction to measures of climatic novelty, by partitioning the climate novelty  
490 according to the type of range shifts from fig. 6B. The study area shows higher novelty in parts of  
491 the range that are currently predicted to be habitable by the species; nevertheless, this does not  
492 translate to an association between types of prediction transition and the distribution of novelty  
493 within the regions undergoing this transition. In other words, the projected uncertainty under  
494 conformal prediction contributes different information when compared to measures of climatic  
495 novelty; specifically, it conveys the uncertainty tied to the model itself.



496 Figure 7: Climate novelty measured as Euclidean distance to the closest contemporary analogue (A); note  
 497 that the scale is square-root transformed, as most areas show low novelty. Distribution of novelty values  
 498 split by the expected transition in occupancy (B); colors are as in fig. 6B.

## 499 Conclusion

500 Conformal prediction, like most SDM methods, is not quite delivering a true estimate of the  
 501 probability of presence (Phillips & Elith 2013). Nevertheless, it brings valuable information, in the  
 502 form of a quantified measure of whether a prediction comes with uncertainty (are both presence  
 503 and absence in the credible set?) in a way that is directly comparable with the non-conformal  
 504 prediction. “Class overlap”, where both presences and absences are observed under the same  
 505 values of the predictions, decreases the predictive performance of models (Valavi et al. 2021) —  
 506 CP is naturally suited at handling this, by assigning the area where overlap occurs to uncertain  
 507 predictions.

508 Transparent communication of uncertainty, meaning that it is both spatially explicit, quantified,  
 509 and expressed under a risk set by the user, is important: we do not expect a fully trained model to  
 510 always be certain, as some areas are genuinely more difficult to predict. For example, small  
 511 organisms are more inherently stochastic (Soininen et al. 2013) any form of stochastic event will  
 512 drive species distribution even when there is strong environmental signal (Mohd et al. 2016) these  
 513 stochastic events can even manifest in areas that are close to the species’ environmental optimum

514 (Dallas et al. 2020). For these reasons, CP can produce interpretable estimates of uncertainty in  
515 species distribution models, and does not require the adoption of additional modeling tools or  
516 paradigms as it functions on an already trained model.

517 CP contributes to dispel what Messeri & Crockett (2024) called the “illusion of understanding”,  
518 which is often associated with ML models: it generates an understanding of the uncertainty from  
519 observations of a pre-trained model, and expresses this uncertainty both in absolute (is the  
520 “presence” event in the credible set?) and relative (is the point estimate of the score for presence  
521 larger than for absence?) terms. Because this technique is computationally efficient and works on  
522 pre-trained models, it opens up the opportunity for more systematic uncertainty quantification  
523 (Zurell et al. 2020) in SDMs. CP, in short, can deliver the “maps of ignorance” that Rocchini et al.  
524 (2011) argued for: how difficult is it to make a prediction for the range at a given risk level is, in  
525 and of itself, an important information to frame the reliability of the results. Finally, CP can  
526 provide guidance on the feedback loop between SDM training and field validation (Johnson et al.  
527 2023) — areas where the range is certain are a much lower priority for sampling.

## 528 Bibliography

- 529 Aldous DJ. 1985. Exchangeability and related topics. In *Lecture Notes in Mathematics*, pp. 1–198.  
530 Berlin, Heidelberg: Springer Berlin Heidelberg
- 531 Allouche O, Tsoar A, Kadmon R. 2006. Assessing the accuracy of species distribution models:  
532 prevalence, kappa and the true skill statistic (TSS). *The journal of applied ecology*. 43(6):1223–  
533 32
- 534 Angelopoulos AN, Bates S. 2023. *Conformal Prediction: A Gentle Introduction*. Hanover, MD: now
- 535 Balayla J. 2020. Prevalence threshold ( $\phi_e$ ) and the geometry of screening curves. *PloS one*.  
536 15(10):e240215

- 537 Barbet-Massin M, Jiguet F, Albert CH, Thuiller W. 2012. Selecting pseudo-absences for species  
538 distribution models: how, where and how many?: How to use pseudo-absences in niche  
539 modelling?. *Methods in ecology and evolution*. 3(2):327–38
- 540 Beale CM, Lennon JJ. 2012. Incorporating uncertainty in predictive species distribution  
541 modelling. *Philosophical transactions of the Royal Society of London. Series B, Biological  
542 sciences*. 367(1586):247–58
- 543 Beery S, Cole E, Parker J, Perona P, Winner K. 2021. Species Distribution Modeling for machine  
544 learning practitioners: A review. *arXiv [cs.LG]*
- 545 Benkendorf DJ, Schwartz SD, Cutler DR, Hawkins CP. 2023. Correcting for the effects of class  
546 imbalance improves the performance of machine-learning based species distribution models.  
547 *Ecological modelling*. 483(110414):110414
- 548 Bylander T. 2002. Estimating Generalization Error on Two-Class Datasets Using Out-of-Bag  
549 Estimates. *Machine learning*. 48(1/3):287–97
- 550 Carlson CJ. 2020. embarcadero: Species distribution modelling with Bayesian additive regression  
551 trees in r. *Methods in ecology and evolution*. 11(7):850–58
- 552 Chen X, Dimitrov NB, Meyers LA. 2019. Uncertainty analysis of species distribution models. *PloS  
553 one*. 14(5):e214190
- 554 Chicco D, Jurman G. 2023. The Matthews correlation coefficient (MCC) should replace the ROC  
555 AUC as the standard metric for assessing binary classification. *BioData mining*. 16(1):4
- 556 Dallas TA, Santini L, Decker R, Hastings A. 2020. Weighing the Evidence for the Abundant-  
557 Center Hypothesis. *Biodiversity informatics*. 15(3):81–91
- 558 Davies SC, Thompson PL, Gomez C, Nephin J, Knudby A, et al. 2023. Addressing uncertainty  
559 when projecting marine species' distributions under climate change. *Ecography*. 2023(11):
- 560 Drake JM. 2014. Ensemble algorithms for ecological niche modeling from presence-background  
561 and presence-only data. *Ecosphere (Washington, D.C.)*. 5(6):1–16

- 562 Elith J. 2019. Species Distribution Modeling
- 563 Fannjiang C, Bates S, Angelopoulos AN, Listgarten J, Jordan MI. 2022. Conformal prediction  
564 under feedback covariate shift for biomolecular design. *Proceedings of the National Academy  
565 of Sciences of the United States of America*. 119(43):e2204569119
- 566 Feng X, Park DS, Walker C, Peterson AT, Merow C, Papeş M. 2019. A checklist for maximizing  
567 reproducibility of ecological niche models. *Nature ecology & evolution*. 3(10):1382–95
- 568 Fick SE, Hijmans RJ. 2017. WorldClim 2: new 1-km spatial resolution climate surfaces for global  
569 land areas: NEW CLIMATE SURFACES FOR GLOBAL LAND AREAS. *International journal  
570 of climatology: a journal of the Royal Meteorological Society*. 37(12):4302–15
- 571 Fitzpatrick MC, Blois JL, Williams JW, Nieto-Lugilde D, Maguire KC, Lorenz DJ. 2018. How will  
572 climate novelty influence ecological forecasts? Using the Quaternary to assess future  
573 reliability. *Global Change Biology*. 24(8):3575–86
- 574 Fontana M, Zeni G, Vantini S. 2020. Conformal Prediction: A unified review of theory and new  
575 challenges. *arXiv [cs.LG]*
- 576 Foxon F. 2024. Bigfoot: If it's there, could it be a bear?. *Journal of zoology (London, England: 1987)*
- 577 Franklin J. 2023. Species distribution modelling supports the study of past, present and future  
578 biogeographies. *Journal of biogeography*. 50(9):1533–45
- 579 Gammerman A, Vovk V, Vapnik V. 1998. Learning by transduction. *Proceedings of the Fourteenth  
580 Conference on Uncertainty in Artificial Intelligence*. 148–55. San Francisco, CA, USA: Morgan  
581 Kaufmann Publishers Inc.
- 582 Janitzka S, Hornung R. 2018. On the overestimation of random forest's out-of-bag error. *PloS one*.  
583 13(8):e201904
- 584 Johnson S, Molano-Flores B, Zaya D. 2023. Field validation as a tool for mitigating uncertainty in  
585 species distribution modeling for conservation planning. *Conservation science and practice*.  
586 5(8):e12978

- 587 Kerr MR, Ordonez A, Riede F, Atkinson J, Pearce EA, et al. 2025. Widespread ecological novelty  
588 across the terrestrial biosphere. *Nature ecology & evolution*. 1–10
- 589 Lei J, Wasserman L. 2013. Distribution-free Prediction Bands for Non-parametric Regression.  
590 *Journal of the Royal Statistical Society. Series B, Statistical methodology*. 76(1):71–96
- 591 Liu C, Newell G, White M. 2016. On the selection of thresholds for predicting species occurrence  
592 with presence-only data. *Ecology and evolution*. 6(1):337–48
- 593 Liu C, White M, Newell G. 2013. Selecting thresholds for the prediction of species occurrence  
594 with presence-only data. *Journal of biogeography*. 40(4):778–89
- 595 Lozier JD, Aniello P, Hickerson MJ. 2009. Predicting the distribution of Sasquatch in western  
596 North America: anything goes with ecological niche modelling. *Journal of biogeography*.  
597 36(9):1623–27
- 598 Mahony CR, Cannon AJ, Wang T, Aitken SN. 2017. A closer look at novel climates: new methods  
599 and insights at continental to landscape scales. *Global change biology*. 23(9):3934–55
- 600 Mesgaran MB, Cousens RD, Webber BL. 2014. Here be dragons: a tool for quantifying novelty due  
601 to covariate range and correlation change when projecting species distribution models.  
602 *Diversity & distributions*. 20(10):1147–59
- 603 Messeri L, Crockett MJ. 2024. Artificial intelligence and illusions of understanding in scientific  
604 research. *Nature*. 627(8002):49–58
- 605 Mohd MH, Murray R, Plank MJ, Godsoe W. 2016. Effects of dispersal and stochasticity on the  
606 presence–absence of multiple species. *Ecological modelling*. 342:49–59
- 607 Neyman J. 1937. Outline of a theory of statistical estimation based on the classical theory of  
608 probability. *Philosophical transactions of the Royal Society of London*. 236(767):333–80
- 609 Oliveira RI, Orenstein P, Ramos T, Romano JV. 2024. Split conformal prediction and non-  
610 exchangeable data. *Journal of machine learning research: JMLR*. 25(225):1–38

- 611 Ordonez A, Riede F, Normand S, Svenning J-C. 2024. Towards a novel biosphere in 2300: rapid  
612 and extensive global and biome-wide climatic novelty in the Anthropocene. *Philosophical*  
613 *transactions of the Royal Society of London. Series B, Biological sciences.* 379(1902):
- 614 Parker EJ, Weiskopf SR, Oliver RY, Rubenstein MA, Jetz W. 2024. Insufficient and biased  
615 representation of species geographic responses to climate change. *Global change biology.*  
616 30(7):e17408
- 617 Petitpierre B, Broennimann O, Kueffer C, Daehler C, Guisan A. 2016. Selecting predictors to  
618 maximize the transferability of species distribution models: lessons from cross-continental  
619 plant invasions. *Global ecology and biogeography: a journal of macroecology.* 26(3):275–87
- 620 Petropoulos F, Hyndman RJ, Bergmeir C. 2018. Exploring the sources of uncertainty: Why does  
621 bagging for time series forecasting work?. *European journal of operational research.*  
622 268(2):545–54
- 623 Phillips SJ, Elith J. 2013. On estimating probability of presence from use-availability or presence-  
624 background data. *Ecology.* 94(6):1409–19
- 625 Poisot T, Bussières-Fournel A, Dansereau G, Catchen MD. 2025. A Julia toolkit for species  
626 distribution data. *EcoEvoRxiv*
- 627 Prescott VA, Marte J, Keller RP. 2025. Performance of alternative methods for generating species  
628 distribution models for invasive species in the Laurentian Great Lakes. *Fisheries.* vuaf12
- 629 Pěknicová J, Berchová-Bímová K. 2016. Application of species distribution models for protected  
630 areas threatened by invasive plants. *Journal for nature conservation.* 34:1–7
- 631 Rocchini D, Hortal J, Lengyel S, Lobo JM, Jiménez-Valverde A, et al. 2011. Accounting for  
632 uncertainty when mapping species distributions: The need for maps of ignorance. *Progress in*  
633 *physical geography.* 35(2):211–26
- 634 Romano Y, Patterson E, Candès E. 2019. Conformalized Quantile Regression. *Neural Information*  
635 *Processing Systems.* 32:3538–48

- 636 Romano Y, Sesia M, Candès EJ. 2020. Classification with valid and adaptive coverage. *Proceedings*  
637       *of the 34th International Conference on Neural Information Processing Systems*. 3581–91. Red  
638 Hook, NY, USA: Curran Associates Inc.
- 639 Roth AE. 1988. Introduction to the Shapley value. In , pp. 1–28. Cambridge University Press
- 640 Sadinle M, Lei J, Wasserman L. 2018. Least Ambiguous Set-Valued Classifiers With Bounded  
641       Error Levels. *Journal of the American Statistical Association*. 114(525):223–34
- 642 Shafer G, Vovk V. 2007. A tutorial on conformal prediction. *Journal of machine learning research:*  
643       *JMLR*. (12):371–421
- 644 Smith JR, Levine JM. 2025. Linking relative suitability to probability of occurrence in presence-  
645       only species distribution models: Implications for global change projections. *Methods in*  
646       *ecology and evolution*
- 647 Soininen J, Korhonen JJ, Luoto M. 2013. Stochastic species distributions are driven by organism  
648       size. *Ecology*. 94(3):660–70
- 649 Soley-Guardia M, Alvarado-Serrano DF, Anderson RP. 2024. Top ten hazards to avoid when  
650       modeling species distributions: a didactic guide of assumptions, problems, and  
651       recommendations. *Ecography*. 2024(4):
- 652 Szeghalmy S, Fazekas A. 2023. A comparative study of the use of stratified cross-validation and  
653       distribution-balanced stratified cross-validation in imbalanced learning. *Sensors (Basel,*  
654       *Switzerland)*. 23(4):
- 655 Thuiller W, Guéguen M, Renaud J, Karger DN, Zimmermann NE. 2019. Uncertainty in ensembles  
656       of global biodiversity scenarios. *Nature communications*. 10(1):1446
- 657 Tibshirani RJ, Barber RF, Candes EJ, Ramdas A. 2019. Conformal Prediction Under Covariate  
658       Shift. *arXiv [stat.ME]*
- 659 Touati S, Radjef MS, Sais L. 2021. A Bayesian Monte Carlo method for computing the Shapley  
660       value: Application to weighted voting and bin packing games. *Computers & operations*  
661       *research*. 125:105094

- 662 Valavi R, Elith J, Lahoz-Monfort JJ, Guillera-Arroita G. 2021. Modelling species presence-only  
663 data with random forests. *Ecography*. 44(12):1731–42
- 664 Vovk V, Nouretdinov I, Manokhin V, Gammerman A. 2018. Cross-conformal predictive  
665 distributions. *Proceedings of the Seventh Workshop on Conformal and Probabilistic Prediction  
666 and Applications*. 91:37–51
- 667 Wadoux AMJ-C, Saby NPA, Martin MP. 2023. Shapley values reveal the drivers of soil organic  
668 carbon stock prediction. *SOIL*. 9(1):21–38
- 669 Zurell D, Elith J, Schröder B. 2012. Predicting to new environments: tools for visualizing model  
670 behaviour and impacts on mapped distributions. *Diversity & distributions*. 18(6):628–34
- 671 Zurell D, Franklin J, König C, Bouchet PJ, Dormann CF, et al. 2020. A standard protocol for  
672 reporting species distribution models. *Ecography*. 43(9):1261–77