

## Article (refereed) - postprint

---

This is the peer reviewed version of the following article:

Chapman, Daniel; Pescott, Oliver L.; Roy, Helen E.; Tanner, Rob. 2019. **Improving species distribution models for invasive non-native species with biologically informed pseudo-absence selection.** *Journal of Biogeography*, 46 (5). 1029-1040, which has been published in final form at <https://doi.org/10.1111/jbi.13555>

This article may be used for non-commercial purposes in accordance with Wiley Terms and Conditions for Use of Self-Archived Versions.

© 2019 John Wiley & Sons Ltd

This version available <http://nora.nerc.ac.uk/522984/>

NERC has developed NORA to enable users to access research outputs wholly or partially funded by NERC. Copyright and other rights for material on this site are retained by the rights owners. Users should read the terms and conditions of use of this material at <http://nora.nerc.ac.uk/policies.html#access>

**This document is the authors' final manuscript version of the journal article, incorporating any revisions agreed during the peer review process. There may be differences between this and the publisher's version. You are advised to consult the publisher's version if you wish to cite from this article.**

The definitive version is available at <http://onlinelibrary.wiley.com/>

Contact CEH NORA team at  
[noraceh@ceh.ac.uk](mailto:noraceh@ceh.ac.uk)

1  
2  
3 1 **Title: Improving species distribution models for invasive non-native species with biologically-**  
4  
5 2 **informed pseudo-absence selection**  
6  
7

8 3 **Running title:** Invasive species distribution models  
9

10 4 **Authors:** Daniel Chapman<sup>1,2</sup>, Oliver L. Pescott<sup>3</sup>, Helen E. Roy<sup>3</sup>, Rob Tanner<sup>4</sup>  
11  
12

13 5 **Institutional affiliations:**  
14

15  
16 6 1 UKRI Centre for Ecology & Hydrology, Edinburgh EH26 0QB, UK  
17  
18

19 7 2 Biological and Environmental Sciences, University of Stirling, Stirling FK9 4LA, UK  
20  
21

22 8 3 UKRI Centre for Ecology & Hydrology, Wallingford OX10 8BB, UK  
23  
24

25 9 4 European and Mediterranean Plant Protection Organisation, 21 Boulevard Richard Lenoir, 75011  
26  
27

28  
29  
30 10 Paris, France

31 11 **Corresponding author:** Daniel Chapman  
32

33 12 **Email addresses:** Daniel Chapman [daniel.chapman@stir.ac.uk](mailto:daniel.chapman@stir.ac.uk), Oliver L. Pescott [olipes@ceh.ac.uk](mailto:olipes@ceh.ac.uk),  
34

35 13 Helen Roy [hele@ceh.ac.uk](mailto:hele@ceh.ac.uk), Rob Tanner [rt@eppo.int](mailto:rt@eppo.int)  
36

37 14 **Acknowledgements:** This research was funded by European Union Life Programme Preparatory  
38  
39 15 project LIFE15 PRE/FR/000001. We thank the Expert Working Groups who performed EPPO Pest  
40  
41 16 Risk Analyses for the five study species and provided us with data and species information to build our  
42  
43 17 models.  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 18 **Abstract**  
4  
5

6 19 **Aim:** We present a novel strategy for species distribution models (SDMs) aimed at predicting the  
7  
8 20 potential distributions of range-expanding invasive non-native species (INNS). The strategy combines  
9  
10 21 two established perspectives on defining the background region for sampling ‘pseudo-absences’ that  
11  
12 22 have hitherto only been applied separately. These are the accessible area, which accounts for dispersal  
13  
14 23 constraints, and the area outside the environmental range of the species and therefore assumed to be  
15  
16 24 unsuitable for the species. We tested an approach to combine these by fitting SDMs using background  
17  
18 25 samples (pseudo-absences) from both types of background.  
19  
20

21 26 **Location:** Global  
22  
23

24 27 **Taxon:** Invasive non-native plants: *Humulus scandens*, *Lygodium japonicum*, *Lespedeza cuneata*,  
25  
26 28 *Triadica sebifera*, *Cinnamomum camphora*  
27  
28

29 29 **Methods:** Presence-background (or presence-only) SDMs were developed for the potential global  
30  
31 30 distributions of five plant species native to Asia, invasive elsewhere and prioritised for risk assessment  
32  
33 31 as emerging INNS in Europe. We compared models where the pseudo-absences were selected from the  
34  
35 32 accessible background, the unsuitable background (defined using biological knowledge of the species’  
36  
37 33 key limiting factors) or from both types of background.  
38  
39

40 34 **Results:** Combining the unsuitable and accessible backgrounds expanded the range of environments  
41  
42 35 available for model fitting and caused biological knowledge about ecological unsuitability to influence  
43  
44 36 the fitted species-environment relationships. This improved the realism and accuracy of distribution  
45  
46 37 projections globally and, generally, within the species’ ranges.  
47  
48

49 38 **Main conclusions:** Correlative SDMs remain valuable for INNS risk mapping and management, but  
50  
51 39 are often criticised for a lack of biological underpinning. Our approach partly addresses this concern by  
52  
53 40 using prior knowledge of species’ requirements or tolerances to define the unsuitable background for  
54  
55 41 modelling, while also accommodating dispersal constraints through considerations of accessibility. It  
56  
57 42 can be implemented with current SDM software and results in more accurate and realistic distribution  
58  
59 43 projections. As such, wider adoption has potential to improve SDMs that support INNS risk assessment.  
60

1  
2  
3 44 **Keywords:** Biomod; climate envelope; ecological niche model; invasive alien species; Maxent; pest  
4  
5 45 risk assessment; presence-absence; presence-only; presence-background; pseudo-absence.  
6  
7  
8 46

9  
10  
11 47 **Introduction**  
12

13 48 Human transport of species beyond their native ranges, leading to biological invasions, is an important  
14  
15 49 driver of ecological change, impacting biodiversity and ecosystem function (Vilà et al., 2011). Decision  
16  
17 50 making about the control and management of invasive non-native species (INNS) is often underpinned  
18  
19 51 by scientific risk assessments, and species distribution models (SDM) are increasingly seen as a  
20  
21 52 valuable tool for this (Jeschke & Strayer, 2008; Václavík & Meentemeyer, 2009; Jiménez-Valverde et  
22  
23 53 al., 2011). The purpose of SDMs applied in this context is to generate risk maps that predict the potential  
24  
25 54 distribution of an INNS as a function of climate and other environmental gradients (Jiménez-Valverde  
26  
27 55 et al., 2011). Specifically, these represent the relative likelihood of establishment should the species be  
28  
29 56 introduced or disperse to each location in the modelled landscape (Elith, 2013). Risk maps can be used  
30  
31 57 for prioritisation of surveillance and management (Peterson & Robins, 2003; Gormley et al., 2011), to  
32  
33 58 estimate the potential spread of emerging INNS in current and future climates (Jiménez-Valverde et al.,  
34  
35 59 2011; Branquart et al., 2016) and to understand the biological and anthropogenic mechanisms governing  
36  
37 60 invasions (Broennimann et al., 2007; Chapman et al., 2014, 2017; Storkey et al., 2014). Clearly, there  
38  
39 61 is a need for robust and accessible SDM tools and methods to ensure the most accurate possible  
40  
41 62 estimation of the potential distributions of INNS.  
42  
43  
44

45  
46 63 Species prioritised for risk assessment in one area have typically already established invasive non-native  
47  
48 64 distributions in other parts of the world (Roy et al., 2014; Branquart et al., 2016; Tanner et al., 2017)  
49  
50 65 necessitating global-scale models and the pooling of distribution data from native and already-invaded  
51  
52 66 ranges (Broennimann & Guisan, 2008; Mainali et al., 2015). Unfortunately species' distributions are  
53  
54 67 rarely documented comprehensively at the spatial resolutions of SDMs (Boakes et al., 2010). Therefore,  
55  
56 68 global-scale models are typically developed using statistical algorithms that contrast the environmental  
57  
58 69 conditions where the species is known to occur with those at 'pseudo-absence' locations sampled from  
59  
60

1  
2  
3 70 a background domain specified by the modeller. Such SDMs are often referred to as presence-only  
4  
5 71 models (Pearce & Boyce, 2006) but we use the term presence-background to differentiate them from  
6  
7 72 ‘one-case’ or true presence-only models that use only the species presences and not the background  
8  
9 73 (Guillera-Arroita et al., 2015). We also differentiate the ‘pseudo-absence’-based presence-background  
10  
11 74 models that are the focus of this study from point process models for species distributions (Warton &  
12  
13 75 Shepherd, 2010). Point process models generalise presence-background models on a more formal  
14  
15 76 statistical basis. However, to our knowledge they are not suitable for grid cell-resolution distribution  
16  
17 77 data, have not been applied for global-scale modelling of INNS and are far less commonly used than  
18  
19 78 well-known presence-background models such as Maxent (Phillips et al., 2008) or the regression and  
20  
21 79 machine learning approaches implemented through software platforms such as Biomod (Thuiller et al.,  
22  
23 80 2009, 2016).

24  
25  
26  
27 81 One important issue when fitting presence-background models to INNS distribution data is that their  
28  
29 82 global distributions are by definition in a non-equilibrium state and are structured by both the species’  
30  
31 83 environmental tolerances and natural and anthropogenic dispersal constraints (Václavík &  
32  
33 84 Meentemeyer, 2009; Elith et al., 2010; Gallien et al., 2010; Chapman et al., 2016). As a consequence,  
34  
35 85 there are suitable but unoccupied regions in which climatic and environmental conditions would permit  
36  
37 86 establishment by the species, but where invasion has not been realised through dispersal. If such regions  
38  
39 87 are included in the background domain, then the model will conflate lack of presence of the species due  
40  
41 88 to dispersal constraints with a lack of presence due to environmental unsuitability, potentially biasing  
42  
43 89 the species-environment relationships and the prediction of potential distributions. Current approaches  
44  
45 90 to reduce this bias emphasise restricting the background domain to an ‘accessible area’ within dispersal  
46  
47 91 range of the occurrences (Elith et al., 2010; Barve et al., 2011; Elith, 2013; Mainali et al., 2015).  
48  
49 92 Although likely to lessen dispersal biases in presence-background models, we suggest this may be  
50  
51 93 overly restrictive for modelling aimed at risk mapping. If background samples are only drawn in close  
52  
53 94 proximity to the occurrences then the range of environmental conditions used to train the model may  
54  
55 95 be insufficient to fully characterise species-environment relationships, impeding the transfer of  
56  
57 96 predictions into other regions (Thuiller et al., 2004; Fitzpatrick & Hargrove, 2009).  
58  
59  
60

1  
2  
3 97 Here, we propose a biologically-informed approach to improve presence-background models for highly  
4  
5 98 dispersal-limited species, such as those undergoing invasive range expansion. The goal is to exclude  
6  
7 99 suitable but unoccupied regions while also maximising the range of environmental conditions used to  
8  
9 100 train the model as well as prior biological knowledge about niche responses to environmental factors.  
10  
11 101 The approach is based on combining two familiar types of background domain – an accessible  
12  
13 102 background in proximity to species’ occurrences (Barve et al., 2011; Mainali et al., 2015) and an  
14  
15 103 unsuitable background outside the environmental envelope of the species (Thuiller et al., 2004;  
16  
17 104 Chefaoui & Lobo, 2007; Le Maitre et al., 2008). Those previous studies have tested both types of  
18  
19 105 background in isolation, but the novel contributions of this study are to combine both types of  
20  
21 106 background, and to emphasise the definition of the unsuitable background using biological knowledge  
22  
23 107 of key limiting factors for the species, e.g. places that do not reach minimum growing temperatures or  
24  
25 108 exceed maximum drought tolerance. By modelling the global distributions of five invasive non-native  
26  
27 109 plants we demonstrate that this constrains the presence-background models to fit more biologically  
28  
29 110 plausible response functions and increases the accuracy of distribution projections.  
30  
31  
32  
33

## 34 111 **Methods**

### 35 112 *Overview*

36  
37 113 Our aim was to compare global-scale presence-background SDMs for INNS developed using  
38  
39 114 background domains defined as only the accessible region, only the unsuitable region, or through our  
40  
41 115 proposed new approach of combining accessible and unsuitable background regions (Figure 1-2).  
42  
43 116 Models were developed to predict the potential distributions of five plant species that are native to  
44  
45 117 temperate and tropical east Asia, highly invasive in other parts of the world and have been prioritised  
46  
47 118 for risk assessment as potentially-emerging invasive non-native plant species in Europe (Branquart et  
48  
49 119 al., 2016; Tanner et al., 2017). The species represent a range of life histories including an annual  
50  
51 120 climbing vine (*Humulus scandens*), a perennial climbing fern (*Lygodium japonicum*), a perennial semi-  
52  
53 121 woody forb (*Lespedeza cuneata*), a deciduous tree (*Triadica sebifera*) and an evergreen tree  
54  
55 122 (*Cinnamomum camphora*).  
56  
57  
58  
59  
60

1  
2  
3 123 *Data for modelling*  
4  
5

6 124 Species occurrences were obtained from a range of sources including Global Biodiversity Information  
7  
8 125 Facility (GBIF), USGS Biodiversity Information Serving Our Nation (BISON), Integrated Digitized  
9  
10 126 Biocollections (iDigBio), iNaturalist, Early Detection and Distribution Mapping System (EDDMapS)  
11  
12 127 and from the members of the European and Mediterranean Plant Protection Organisation (EPPO) expert  
13  
14 128 working groups conducting Pest Risk Analyses for the region. With these experts, we scrutinised  
15  
16 129 occurrence records and removed any that appeared dubious, casual or cultivated (e.g. botanic gardens)  
17  
18 130 or where the georeferencing was too imprecise (e.g. country or island centroids). The remaining records  
19  
20 131 were gridded at a 0.25 x 0.25 degree resolution for global modelling and randomly partitioned into  
21  
22 132 training and testing datasets comprising 80% and 20% of the grid cells, respectively. As a proxy for  
23  
24 133 plant recording effort, the total number of vascular plant records (phylum Tracheophyta) per grid cell  
25  
26 134 was also obtained from GBIF (see Appendix S1 in Supporting Information).

27  
28  
29 135 Three predictor variables, derived from WorldClim v1.4 (Hijmans et al., 2005), were selected to  
30  
31 136 represent basic constraints on plant distributions. These were mean temperature of the warmest quarter  
32  
33 137 (Bio10, °C) reflecting the growing season thermal regime, mean minimum temperature of the coldest  
34  
35 138 month (Bio6, °C) reflecting exposure to winter cold and the climatic moisture index (CMI, ratio of  
36  
37 139 annual precipitation, Bio12, to potential evapotranspiration, then ln + 1 transformed) reflecting drought  
38  
39 140 stress. Potential evapotranspiration was estimated following Zomer et al. (2008).

40  
41  
42  
43 141 *Definition of the background domains*  
44

45  
46 142 Background samples (pseudo-absences) were drawn from two distinct regions – an accessible region  
47  
48 143 and a region considered to be environmentally unsuitable for the species based on knowledge of its  
49  
50 144 tolerances or requirements (Figures 1 and 2). Though both types of background represent established  
51  
52 145 concepts within distribution modelling, to our knowledge, this is the first study to test whether  
53  
54 146 modelling is improved by combining both types of background domain.

55  
56  
57 147 The accessible background attempts to cover only the region where the species has had opportunity to  
58  
59 148 disperse and sample the environment (Thuiller et al., 2004; VanDerWal et al., 2009; Barve et al., 2011;  
60

1  
2  
3 149 Mainali et al., 2015). It has generally been defined as a zone around the occurrence data, which could  
4  
5 150 be selected statistically or informed by dispersal abilities of the species (Elith, 2013; Senay et al., 2013).  
6  
7 151 For invasive non-native species, the size of the accessible region will generally be more limited in the  
8  
9 152 invaded range than the native one, assuming stronger dispersal constraints associated with shorter  
10  
11 153 residence time (Mainali et al., 2015). In our application, we defined the native accessible areas using a  
12  
13 154 400 km geodesic buffer around the minimum convex polygon bounding all native occurrences (Figure  
14  
15 155 1a). In the non-native region, we used a conservative 4-cell neighbourhood around each occurrence grid  
16  
17 156 cell, equivalent to a ~30 km buffer (Figure 1b). Though somewhat arbitrary, these buffer sizes are  
18  
19 157 consistent with ones performing well in other presence-background SDM studies (VanDerWal et al.,  
20  
21 158 2009; Mainali et al., 2015) and a sensitivity analysis showed model outputs were not strongly influenced  
22  
23 159 by the choice of native buffer size (see Appendix S5).

24  
25  
26  
27 160 The unsuitable background concept originates from existing ideas about sampling pseudo-absences  
28  
29 161 only outside of the environmental envelope in which species' presences are found (Thuiller et al., 2004;  
30  
31 162 Chefaoui & Lobo, 2007; Le Maitre et al., 2008; Senay et al., 2013). The rationale is to produce training  
32  
33 163 datasets that maximise the distinctiveness of suitable environmental conditions from the background  
34  
35 164 and therefore boost the model discrimination. However, it may also reduce model accuracy within the  
36  
37 165 environmental and geographical range of the species (Acevedo et al., 2012). These previous studies  
38  
39 166 simply screened out the ranges of all environmental variables at presence locations, or used preliminary  
40  
41 167 modelling to determine unsuitable regions. However, in this study we instead used prior biological  
42  
43 168 knowledge and expert opinion about the species' limiting factors to define the unsuitable conditions  
44  
45 169 (Figures 1 and 2) in the expectation that this biological information would be captured in the fitted  
46  
47 170 species-environment relationships. Appropriate rules to define unsuitability were determined in  
48  
49 171 consultation with species experts participating in their EPPO expert working groups. Their expert  
50  
51 172 judgement informed us on the type of limit deemed to be most important for the species in different  
52  
53 173 parts of its range (e.g. summer cold, drought), followed by identification of key thresholds from the  
54  
55 174 literature and comparison with extreme values at the occurrence locations of the species (see Appendix  
56  
57 175 S2).

1  
2  
3 176 *Sampling from the background domains*  
4  
5

6 177 We obtained background samples from both the accessible region and from the unsuitable region  
7  
8 178 outside of the accessible region for each species (Figures 1-2). The effect was therefore to exclude  
9  
10 179 potentially suitable but inaccessible regions from the combined background sample. For each of the  
11  
12 180 five species in this study, ten replicate background samples were generated in order to reduce sampling  
13  
14 181 variation (Barbet-Massin et al., 2012). Presence-background models were developed for each  
15  
16 182 background sample and then their predictions were averaged.  
17  
18

19 183 The accessible region was sampled using target group sampling to reduce bias in the observed  
20  
21 184 distribution due to spatial sampling effort variation (Phillips, 2009; Ranc et al., 2017). This involves  
22  
23 185 weighting the background sampling by the recording density of a broader taxonomic group, which is  
24  
25 186 assumed to represent recording bias for the focal species. In our modelling we used the GBIF record  
26  
27 187 density of vascular plants (Tracheophyta) as a target group to weight background sampling. For  
28  
29 188 evaluating the models by cross-validation, a randomly selected 20% of the accessible area for each  
30  
31 189 species was added to its testing dataset and reserved from background sampling for model fitting. From  
32  
33 190 the remaining accessible area, we drew the same number of background samples as there were  
34  
35 191 occurrences (Barbet-Massin et al., 2012), weighted by the vascular plant record density as a target  
36  
37 192 group. This ensured that the accessible area background sample used for model fitting contained the  
38  
39 193 same degree of recording bias as the occurrence data, assuming the proxy for recording effort was  
40  
41 194 appropriate.  
42  
43  
44

45 195 The unsuitable region was sampled with simple random sampling because we considered that recording  
46  
47 196 bias is not a relevant consideration in environments in which the species cannot occur. In other words,  
48  
49 197 we were confident of absence in the unsuitable regions. Although we could have nevertheless applied  
50  
51 198 target group sampling, random sampling has the potential advantage of accumulating background  
52  
53 199 samples from unsuitable environments where there is little survey effort (e.g. very cold conditions),  
54  
55 200 resulting in the widest range of environments from which to model species-environment relationships.  
56  
57 201 For model fitting, 3000 random samples were taken from the unsuitable region. If the unsuitable region  
58  
59 202 overlapped with the accessible region, accessible parts of the unsuitable region were excluded. A  
60

1  
2  
3 203 sensitivity analysis on the number of unsuitable background samples showed that the number of  
4  
5 204 sampling points was not critical to model performance (see Appendix S5).  
6  
7

8 205 *Ensemble presence-background modelling*  
9

10 206 For each species, presence-background models were developed using background samples from only  
11  
12 207 the accessible area, only the unsuitable area or using the combined background samples from both the  
13  
14 208 accessible and unsuitable areas. In all cases, model performance was evaluated by cross validation,  
15  
16 209 using model predictions for 20% of the occurrences, accessible area and unsuitable area that were  
17  
18 210 reserved from model fitting (the testing dataset).  
19  
20

21  
22 211 Ensemble models were fitted using BIOMOD (biomod2 R package v3.3-7) (Thuiller et al., 2009, 2016)  
23  
24 212 using seven statistical algorithms: generalised linear models (GLM) with linear and quadratic terms for  
25  
26 213 each predictor, generalised additive models (GAM) with a maximum of four degrees of freedom per  
27  
28 214 variable, multivariate adaptive regression splines (MARS), generalised boosting models (GBM),  
29  
30 215 random forests (RF), artificial neural networks (ANN) and Maxent (Phillips et al., 2008). These were  
31  
32 216 combined into an ensemble model by scaling their predictions with a binomial GLM and then averaging  
33  
34 217 them weighted by predictive AUC scores within the training data (80:20% random split). AUC is  
35  
36 218 commonly used for ensemble model weighting and is the BIOMOD default option (Thuiller et al., 2009,  
37  
38 219 2016). Although AUC does not provide an objective measure of model performance for presence-only  
39  
40 220 models (Lobo, 2008) it is informative about the relative discrimination abilities of different algorithms  
41  
42 221 evaluated on the same data. It also provides a conservative model weighting scheme, since a perfect  
43  
44 222 model (AUC=1) will have only twice the weight of a random model (AUC=0.5). Therefore, we ensured  
45  
46 223 poorly performing algorithms did not disproportionately affect the weighted average by rejecting them  
47  
48 224 from the ensemble. Rejection was based on modified z-scores for their predictive AUC (Crosby, 1993)  
49  
50 225 with algorithms with  $z < -1$  being rejected.  
51  
52  
53

54 226 The importance of each variable to model fitting was estimated through the BIOMOD default procedure  
55  
56 227 (Thuiller et al., 2009, 2016). Species-environment relationships were examined by constructing  
57  
58 228 univariate response curves where predictions of the ensemble model were made while fixing the other  
59  
60

1  
2  
3 229 variables at typical suitable values representing the median in the presence grid cells. Global projections  
4  
5 230 of the ensemble models were restricted to where the environmental predictors lay inside the ranges used  
6  
7 231 in model training, avoiding model extrapolation (Fitzpatrick & Hargrove, 2009).  
8  
9

10 232 Models based only on the accessible or the unsuitable background were compared with those based on  
11  
12 233 the combined accessible and unsuitable background in a standardised cross validation. To do this we  
13  
14 234 used calculated the AUC for model projections on the 20% of occurrences versus the 20% of the  
15  
16 235 accessible background that was reserved from model fitting and versus 20% of the accessible and  
17  
18 236 unsuitable background. This allows comparison of projection accuracy within the range of the species  
19  
20 237 and globally. As mentioned above, AUC in this context is informative about the relative discrimination  
21  
22 238 power of different model specifications on the same data.  
23  
24  
25  
26

## 27 239 **Results**

28  
29 240 Adequate numbers of grid cells with presences were obtained for modelling the five study species (695  
30  
31 241 for *Cinnamomum camphora*, 754 for *Humulus scandens*, 1723 for *Lespedeza cuneata*, 975 for  
32  
33 242 *Lygodium japonicum* and 855 for *Triadica sebifera*) (see Appendix S2). In most cases, cross-validated  
34  
35 243 AUC indicated that models trained using samples from the combined accessible and unsuitable  
36  
37 244 background were more accurate than those trained using only the individual accessible or unsuitable  
38  
39 245 backgrounds (Table 1 and see Appendix S3). This was most clearly seen for global projections of the  
40  
41 246 model, where the combined models had the highest AUC values for all five species (Table 1). The  
42  
43 247 probability of the combined background model having the highest AUC of the three model types for all  
44  
45 248 five species by chance is  $P = 0.004$ . For projections within the accessible range of the species, models  
46  
47 249 sampling the combined accessible and unsuitable background were equally or marginally more accurate  
48  
49 250 than models using only the accessible background in four out of five species, and always performed  
50  
51 251 better than models using only the unsuitable background (Table 1).  
52  
53  
54

55 252 Models using only the accessible background spanned a narrower range of suitability values and  
56  
57 253 environmental conditions than the other two background specifications, and therefore their response  
58  
59 254 curves were only constructed over a narrow range and provided a less clear distinction between high  
60

1  
2  
3 255 and low suitability (Figure 3). Models using only the unsuitable background generated response curves  
4  
5 256 that essentially discriminated unsuitable from suitable conditions completely, but did not capture  
6  
7 257 suitability gradients within the range of the species (Figure 3). By contrast, the models based on the  
8  
9 258 combined accessible and unsuitable background yielded partial response curves intermediate in form to  
10  
11 259 the two individual background specifications (Figure 3), using information from the accessible  
12  
13 260 background region to characterise gradients in suitability within the environmentally-suitable range of  
14  
15 261 the species, and using the unsuitable background to identify conditions in which the species very rarely  
16  
17 262 occurs. In some cases the models using combined accessible and unsuitable backgrounds yielded  
18  
19 263 response curves that differed markedly from those of the accessible background models. This was most  
20  
21 264 clearly seen in the responses of *Cinnamomum camphora* and *Lygodium japonicum* to low moisture  
22  
23 265 (CMI), *Lespedeza cuneata* to low winter temperature (Bio6) and *Cinnamomum camphora* to high  
24  
25 266 summer temperature (Bio10) (Figure 3).

26  
27  
28  
29 267 Projections of potential non-native ranges from the models were strongly influenced by the choice of  
30  
31 268 background specification (Figures 4 and 5, see Appendix S4 for global and native range projections).  
32  
33 269 As was seen for the response curves, models based only on the accessible background generally made  
34  
35 270 a gradual delineation between very low and high suitability, while models based on the unsuitable  
36  
37 271 background made very sharp delineations and predicted larger invadable regions. Projections of models  
38  
39 272 using the combined accessible and unsuitable backgrounds were intermediate in form, and represented  
40  
41 273 gradients in suitability within the invaded regions as well as learning from the unsuitable background  
42  
43 274 to rule out occurrence in those regions. For example, in North America the models using only the  
44  
45 275 accessible background predicted suitability for *Cinnamomum camphora*, *Triadica sebifera* and  
46  
47 276 *Lygodium japonicum* invasion in arid parts of south western USA. By contrast, models combining the  
48  
49 277 accessible and unsuitable backgrounds suggested these regions were unsuitable for invasion (Figure 4).  
50  
51 278 Similar effects could be seen in Europe, principally in terms of the effects of unsuitability rules about  
52  
53 279 low winter temperature restricting suitability in central and Eastern Europe and rules about drought  
54  
55 280 sensitivity restricting occurrence in Iberia (Figure 5). The projections also illustrated that models based  
56  
57  
58  
59  
60

1  
2  
3 281 on only the accessible background were affected by extrapolation issues, resulting in gaps in the  
4  
5 282 projected risk maps.  
6  
7  
8

9 283 **Discussion**

10  
11  
12 284 Strategies for selecting background samples or pseudo-absences for presence-background species  
13  
14 285 distribution models have received a great deal of attention (e.g. Thuiller et al., 2004; Chefaoui & Lobo,  
15  
16 286 2007; VanDerWal et al., 2009; Barve et al., 2011; Barbet-Massin et al., 2012). The novel contribution  
17  
18 287 of this study is to combine two different perspectives on defining the background region that have  
19  
20 288 hitherto been considered separately. These perspectives are the accessible area (Barve et al., 2011) and  
21  
22 289 the area outside the environmental range of the species, and therefore assumed to be unsuitable for the  
23  
24 290 species (Thuiller et al., 2004). Previous work on modelling invasive non-native species has generally  
25  
26 291 either emphasised the usefulness of the former for accommodating dispersal constraints (Mainali et al.,  
27  
28 292 2015) or evaluated the latter as a way of boosting the discrimination between suitable and unsuitable  
29  
30 293 habitat (Le Maitre et al., 2008). To our knowledge, the only previous attempt to jointly consider both  
31  
32 294 perspectives did so in a more limited way than this study, by excluding parts of the accessible region  
33  
34 295 that were outside the environmental range of the species (Senay et al., 2013). Here, we tested a new  
35  
36 296 approach in which separate background samples were obtained from the accessible region, regardless  
37  
38 297 of environmental values, and from an unsuitable region defined using prior biological knowledge. By  
39  
40 298 modelling the global distributions of five invasive non-native plant species we conclude that the new  
41  
42 299 strategy performed better for projection of regional and global potential distributions than when models  
43  
44 300 were fitted with just the accessible region or just the unsuitable region.  
45  
46  
47

48  
49 301 This was evidenced by a consistent improvement in cross-validated discrimination power when the  
50  
51 302 modelling sampled from a background combining accessible and biologically-informed unsuitable  
52  
53 303 regions. This was most clearly seen in the global projections, where the combined background models  
54  
55 304 always performed better than models using just the accessible or just the unsuitable background. For  
56  
57 305 projections within the species' accessible range the combined background models gave consistently  
58  
59 306 more accurate projections than models based only on the unsuitable background, and generally  
60

1  
2  
3 307 performed as well as or marginally better than models trained only on the accessible background. Our  
4  
5 308 expectation was that the combined background modelling strategy would not improve discrimination  
6  
7 309 within the range of a species over models trained on the accessible region. Indeed, previous studies have  
8  
9 310 found that large geographical background domains increase the power of SDMs to model species' broad  
10  
11 311 geographic ranges but decrease their representation of suitability gradients within the range (Thuiller et  
12  
13 312 al., 2004; VanDerWal et al., 2009). Unlike previous studies, our approach may have resulted in  
14  
15 313 marginally improved performance for both purposes because we explicitly tried to exclude 'suitable-  
16  
17 314 but-not-reached' locations from the larger background region by restricting it to locations considered  
18  
19 315 environmentally unsuitable. As such, we suggest that biologically-informed specification of a large  
20  
21 316 modelling domain may reduce the trade-off between prediction of suitability gradients at large and  
22  
23 317 small spatial scales. Further testing is required to determine whether a similar strategy would also  
24  
25 318 benefit models for native as well as non-native species distribution models, but in principle our new  
26  
27 319 strategy should confer similar advantages.

30  
31 320 The influence of the accessible and unsuitable backgrounds on species-environment relationships was  
32  
33 321 clearly seen in the response curves and projections of the different models. The combination of  
34  
35 322 unsuitable and accessible backgrounds had four clear effects, when compared to the models using only  
36  
37 323 the accessible background. First, it 'anchored' the curves by constraining the models to fit near-zero  
38  
39 324 suitability where the climate variables exceeded the thresholds of the species, providing a more  
40  
41 325 pronounced delineation of suitability gradients. Second, the response curves spanned a much wider  
42  
43 326 range of environmental conditions than were found in the accessible background, which has previously  
44  
45 327 been shown to be important for accurate spatial and temporal transfer of species distribution models  
46  
47 328 (Guevara et al., 2017). Sampling unsuitable conditions only from within the accessible part of the  
48  
49 329 species range would therefore require a greater amount of model extrapolation than our strategy does.  
50  
51 330 Third, the response curves were less complex or multi-modal than those from models using only the  
52  
53 331 accessible background (see responses for high CMI), which is more consistent with niche theory  
54  
55 332 (Austin, 2002). Fourth, the response curves generally reflected prior assumptions about environmental  
56  
57 333 limitation of the species and as such were more consistent with ecological understanding of the species.  
58  
59  
60

1  
2  
3 334 For instance, combined background models for *Cinnamomum camphora*, *Lygodium japonicum* and  
4  
5 335 *Triadica sebifera* estimated a strong limitation by low moisture availability (CMI), precluding potential  
6  
7 336 establishment in arid regions such as south west USA. These responses were not estimated by the model  
8  
9 337 based only on the accessible background, but are consistent with empirical demonstrations of water  
10  
11 338 stress reducing growth and survival of these species. For example, shoot growth of *C. camphora* is 30%  
12  
13 339 lower at 40% field water capacity than at 80% (Zhao et al., 2006), water restriction suppresses *T.*  
14  
15 340 *sebifera* seedling growth by 30-80% (Barrilleaux & Grace, James, 2000) and its seedlings wilt and die  
16  
17 341 in arid western USA unless planted in moist micro-habitats such as river banks (Bower et al., 2009).  
18  
19 342 Similarly, combining the accessible and unsuitable backgrounds led to models that strongly limited  
20  
21 343 suitability of *Lespedeza cuneata* by very cold winters, consistent with known frost sensitivity of the  
22  
23 344 species especially in relation to late spring frosts (Gucker, 2010). The only case where the response  
24  
25 345 curves did not always follow the rules defining unsuitability was for limitation by extremely high  
26  
27 346 summer temperature. This may be because of a correlation between high summer and winter  
28  
29 347 temperatures, the latter being limiting when high summer temperature was not. This suggests our  
30  
31 348 approach may have sensitivity to collinearity in model predictors that requires further investigation  
32  
33 349 (Dormann et al., 2012).  
34  
35  
36  
37 350 Nevertheless, the broader conclusion is that sampling from an unsuitable background, in addition to an  
38  
39 351 accessible background, forces the statistical models to learn species-environment relationships that  
40  
41 352 reflect the prior knowledge of the species' tolerances or niche requirements used to define the unsuitable  
42  
43 353 domain. As such, our approach offers a simple way of incorporating prior biological knowledge into  
44  
45 354 correlative species distribution models, and as such can address the common criticism that they lack  
46  
47 355 strong biological underpinning (Austin, 2002; Dormann et al., 2011; Chapman et al., 2014). While there  
48  
49 356 are more sophisticated approaches available for doing this using Bayesian models in which prior  
50  
51 357 estimates of niche parameters can be specified (Talluto et al., 2015), a major advantage of the approach  
52  
53 358 developed here is that it is implemented by manipulating the input data to standard distribution model  
54  
55 359 software such as Biomod (Thuiller et al., 2009) or Maxent (Phillips et al., 2008) and all regression and  
56  
57 360 machine learning methods. As such it is simple to implement with techniques that most modellers are  
58  
59  
60

1  
2  
3 361 already familiar with and can quickly be applied in a standard way across species. This is especially  
4  
5 362 useful when risk assessments are being performed across large numbers of invasive non-native species  
6  
7 363 and require consistent judgements about establishment risk (Branquart et al., 2016; Tanner et al., 2017).  
8  
9  
10 364 Sensitivity analyses suggested that our findings were not overly sensitive to the size of the accessible  
11  
12 365 region, number of background samples or precise rules for determining unsuitable conditions (see  
13  
14 366 Appendix S5). We recommend that similar sensitivity analyses are performed when applying our  
15  
16 367 approach to other species, since previous studies have found these factors can strongly influence  
17  
18 368 distribution model performance (Barve et al., 2011; Barbet-Massin et al., 2012). However, success of  
19  
20 369 the modelling approach likely relies on careful selection of the appropriate environmental limits to  
21  
22 370 define the unsuitable region in the modelling (Le Maitre et al., 2008). A strength of this study is that it  
23  
24 371 was done in consultation with experts performing risk assessments for invasion of Europe by the  
25  
26 372 species. These experts were able to provide guidance on the key limiting factors relevant for different  
27  
28 373 parts of the invaded and native ranges of the species. Some of the species have been well studied in  
29  
30 374 their other invaded ranges and we were able to draw upon previous experimental studies that had  
31  
32 375 determined tolerance thresholds for the species (see Appendix S2). Where this information was lacking,  
33  
34 376 we used upper or lower bounds on the environmental values at the species presences to define thresholds  
35  
36 377 for modelling. Even where empirical estimates of threshold values were available, we still recommend  
37  
38 378 checking for consistency with environmental values at the distribution data, since species-environment  
39  
40 379 relationships are highly scale-dependent (Siefert et al., 2012) and species can occupy broadly unsuitable  
41  
42 380 regions if suitable micro-habitats are available. Given the reliance on prior studies or expert judgement  
43  
44 381 about species' limiting factors or tolerances, our methods are probably most suitable for relatively well  
45  
46 382 known species and less applicable to species where knowledge of its environmental limits are lacking.  
47  
48 383 However, regional risk assessments for emerging invasive non-native species generally prioritise  
49  
50 384 species that behave invasively in other parts of the world (Roy et al., 2014; Branquart et al., 2016;  
51  
52 385 Tanner et al., 2017) suggesting that our modelling approach might be widely applicable for species of  
53  
54 386 concern.  
55  
56  
57  
58  
59  
60

1  
2  
3 387 Risk assessment is a critical tool in the management of emerging invasive non-native species and  
4  
5 388 requires robust prediction of where is vulnerable to ongoing species establishment and spread (Keller  
6  
7 389 et al., 2007; Jiménez-Valverde et al., 2011). This study shows that defining the model background to  
8  
9 390 accommodate considerations of accessibility as well as prior biological knowledge of environmental  
10  
11 391 unsuitability has the potential to improve global-scale presence-background models for emerging  
12  
13 392 invasive non-native species. The methods developed and tested here are fully implemented by  
14  
15 393 manipulating the model input data, and as such they can be implemented simply using standard  
16  
17 394 presence-background modelling software. Furthermore, they result in presence-background models that  
18  
19 395 are more strongly underpinned by biological knowledge rather than being solely driven by distribution  
20  
21 396 data, which are often incomplete and biased. As such, wider adoption of these approaches should  
22  
23 397 improve global-scale modelling of invasive non-native species distributions, contributing to more  
24  
25 398 accurate risk assessment and better management of their impacts.  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

399 **Tables**

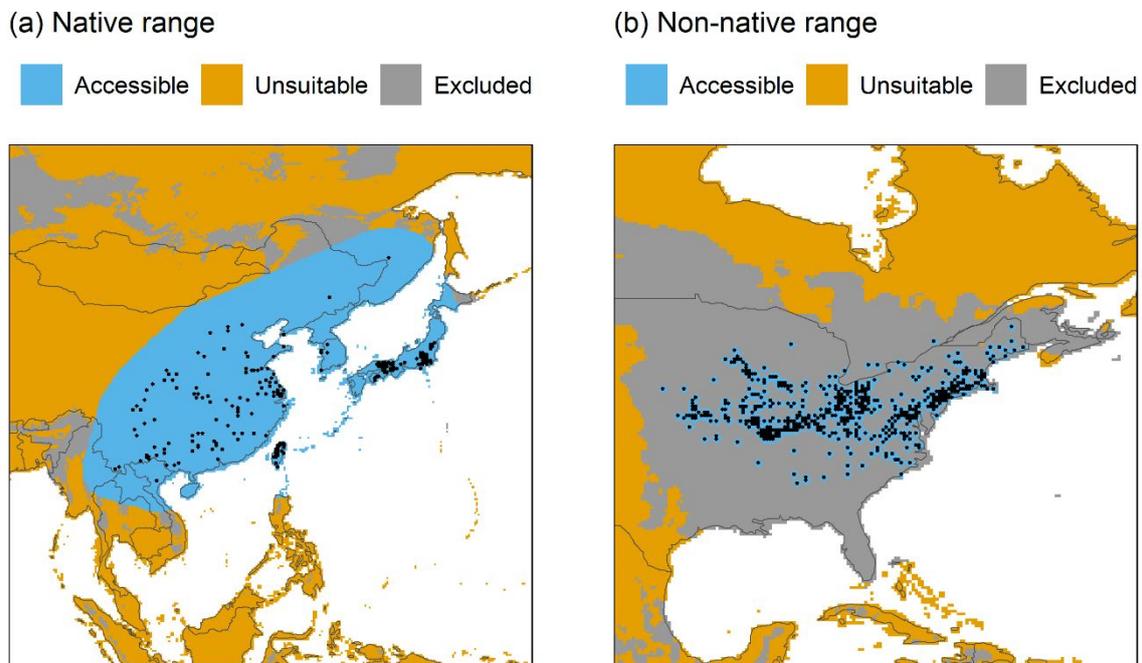
400 **Table 1.** Cross-validated discrimination performance of ensemble model projections for the potential  
 401 global distribution of five plant species developed using different background region specifications (A  
 402 = accessible background, U = unsuitable background, A&U = combined accessible and unsuitable  
 403 background). Discrimination performance is the cross-validated AUC (Area Under the receiver-  
 404 operator Curve) and its standard deviation in parentheses for model predictions on 20% of the  
 405 occurrences, accessible background and unsuitable background that were reserved from model fitting.  
 406 For presence-only data AUC is the probability that a species presence has a higher projected suitability  
 407 than a background sample.

Species	Accuracy in the species range (AUC in accessible background)			Global accuracy (AUC in the accessible and unsuitable backgrounds)		
	A	U	A&U	A	U	A&U
<i>Cinnamomum camphora</i>	0.664 (0.019)	0.581 (0.020)	0.669 (0.020)	0.857 (0.008)	0.981 (0.001)	0.985 (0.001)
<i>Humulus scandens</i>	0.742 (0.020)	0.669 (0.017)	0.737 (0.021)	0.977 (0.004)	0.979 (0.001)	0.982 (0.002)
<i>Lespedeza cuneata</i>	0.899 (0.006)	0.860 (0.006)	0.899 (0.006)	0.979 (0.003)	0.977 (0.002)	0.983 (0.002)
<i>Lygodium japonicum</i>	0.852 (0.013)	0.758 (0.014)	0.852 (0.013)	0.955 (0.007)	0.979 (0.001)	0.987 (0.001)
<i>Triadica sebifera</i>	0.762 (0.017)	0.673 (0.016)	0.777 (0.017)	0.853 (0.007)	0.984 (0.001)	0.989 (0.001)

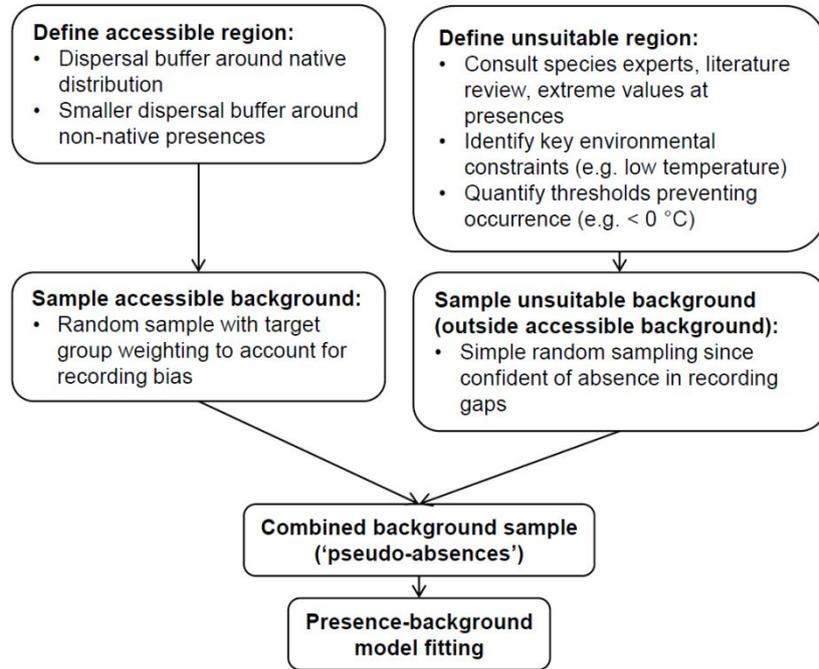
408

1  
2  
3 409 **Figures**  
4  
5

6 410 **Figure 1.** Parts of the regions from which background samples (pseudo-absences) were drawn for  
7  
8 411 modelling one of the five species, *Humulus scandens*. Shading shows the accessible background, where  
9  
10 412 the species is assumed to have had chance to disperse to and sample, and the unsuitable background,  
11  
12 413 defined using biological information on the key limiting factors of the species (see Appendix S2).  
13  
14 414 Potentially suitable, but inaccessible locations were excluded from the modelling (a) The Asian native  
15  
16 415 range of the species, where accessibility was defined with a buffer around the minimum convex polygon  
17  
18 416 of the occurrences. (b) The North American part of the invaded range, where accessibility was highly  
19  
20 417 restricted to represent stronger dispersal constraints during the invasive range expansion.  
21  
22

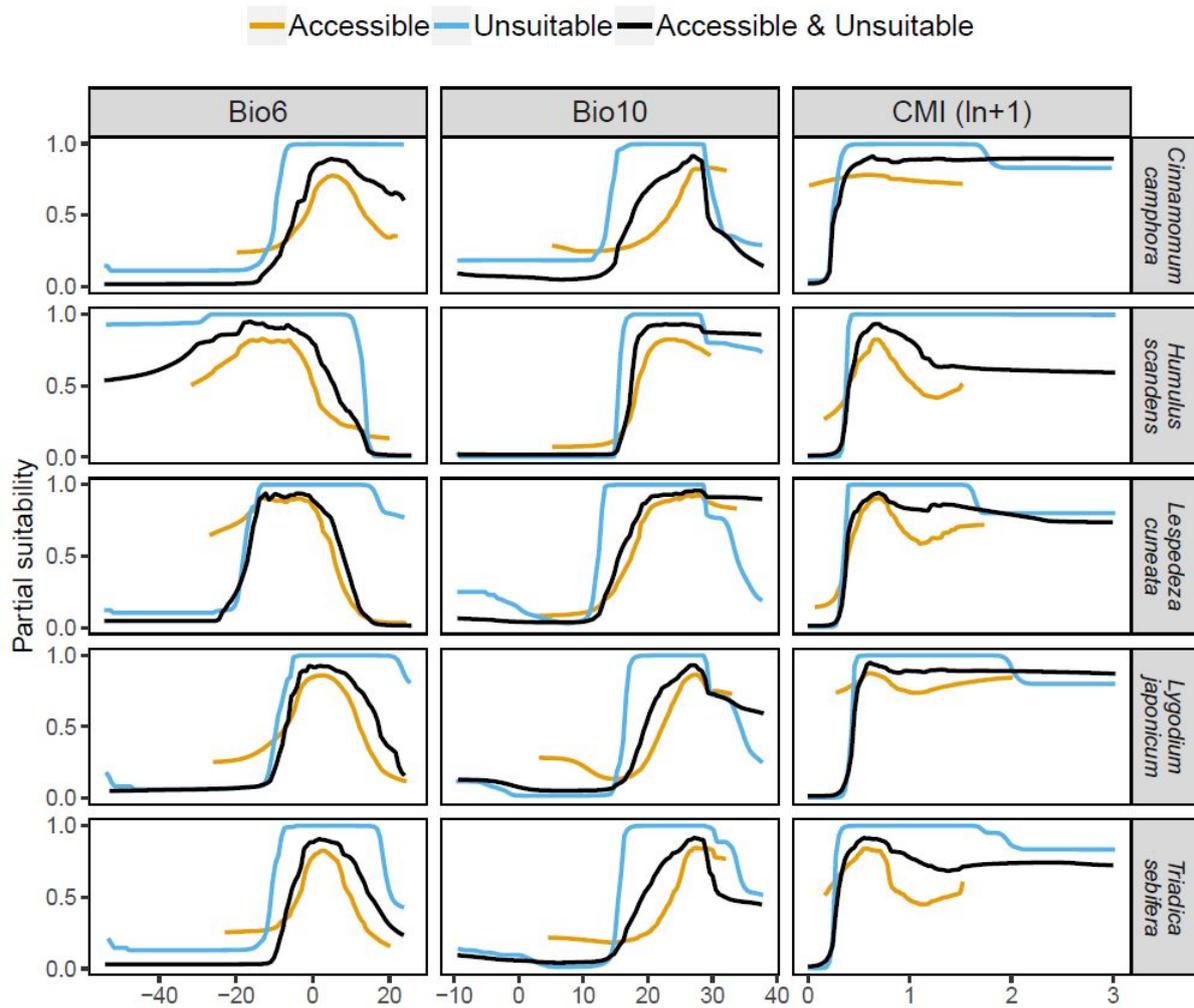


1  
2  
3 419 **Figure 2.** Flow chart for implementing the biologically-informed pseudo-absence selection for  
4  
5 420 presence-background modelling of invasive non-native species.  
6  
7  
8

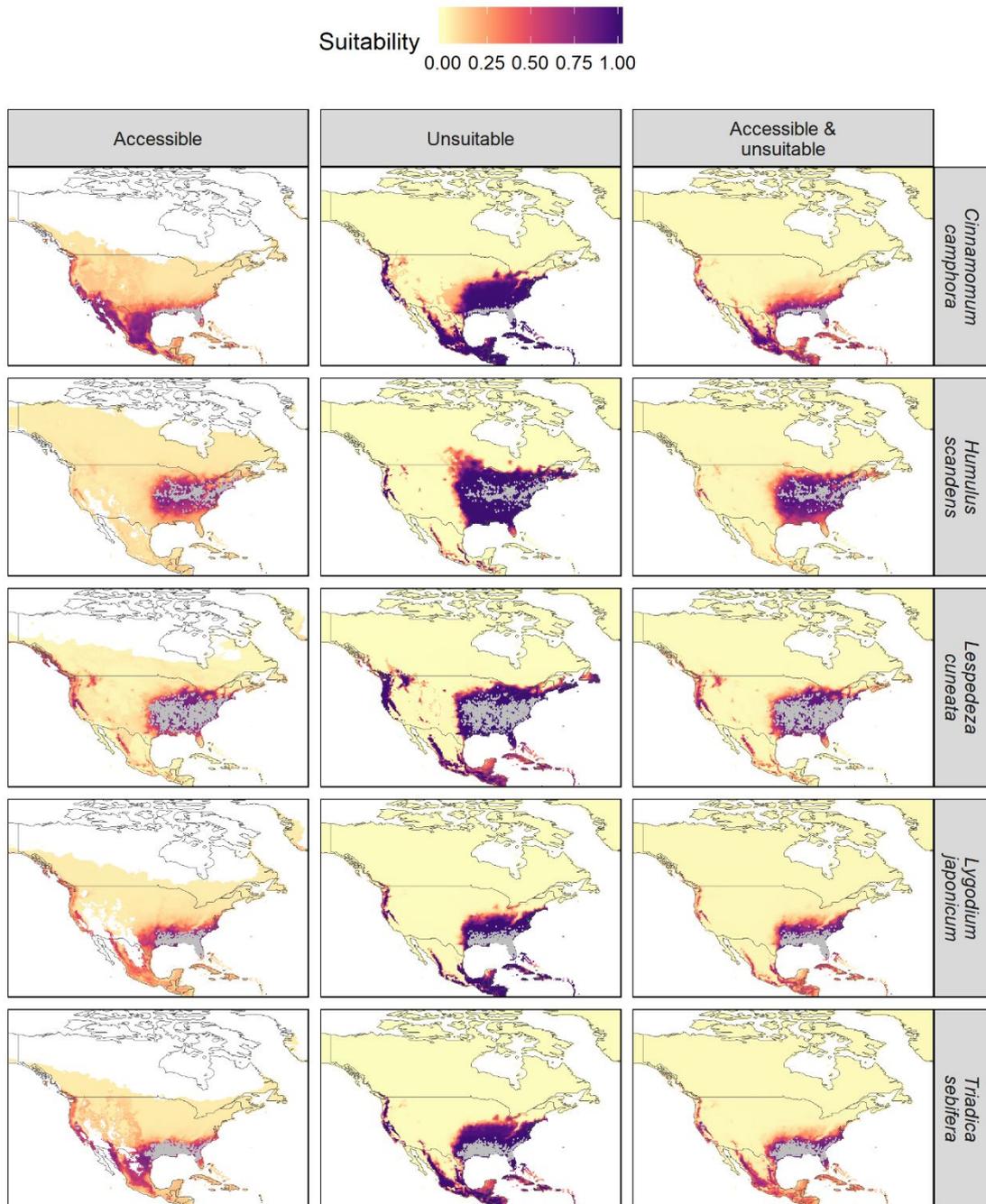


29 421

1  
 2  
 3 422 **Figure 3.** Partial response plots fitted by the ensemble models showing the predicted suitability when  
 4  
 5 423 other variables are fixed at suitable values for the species (medians in the presence grid cells). Curves  
 6  
 7 424 span the range of the variables in the training data. Curve colour differentiates the models with  
 8  
 9 425 background domains based only on the accessible region and those including the unsuitable region.  
 10  
 11 426 Variable codes: Bio6 = mean minimum temperature of the coldest month (°C); Bio10 = mean  
 12  
 13 427 temperature of the warmest quarter (°C); CMI = climatic moisture index (ln+1 transformed).

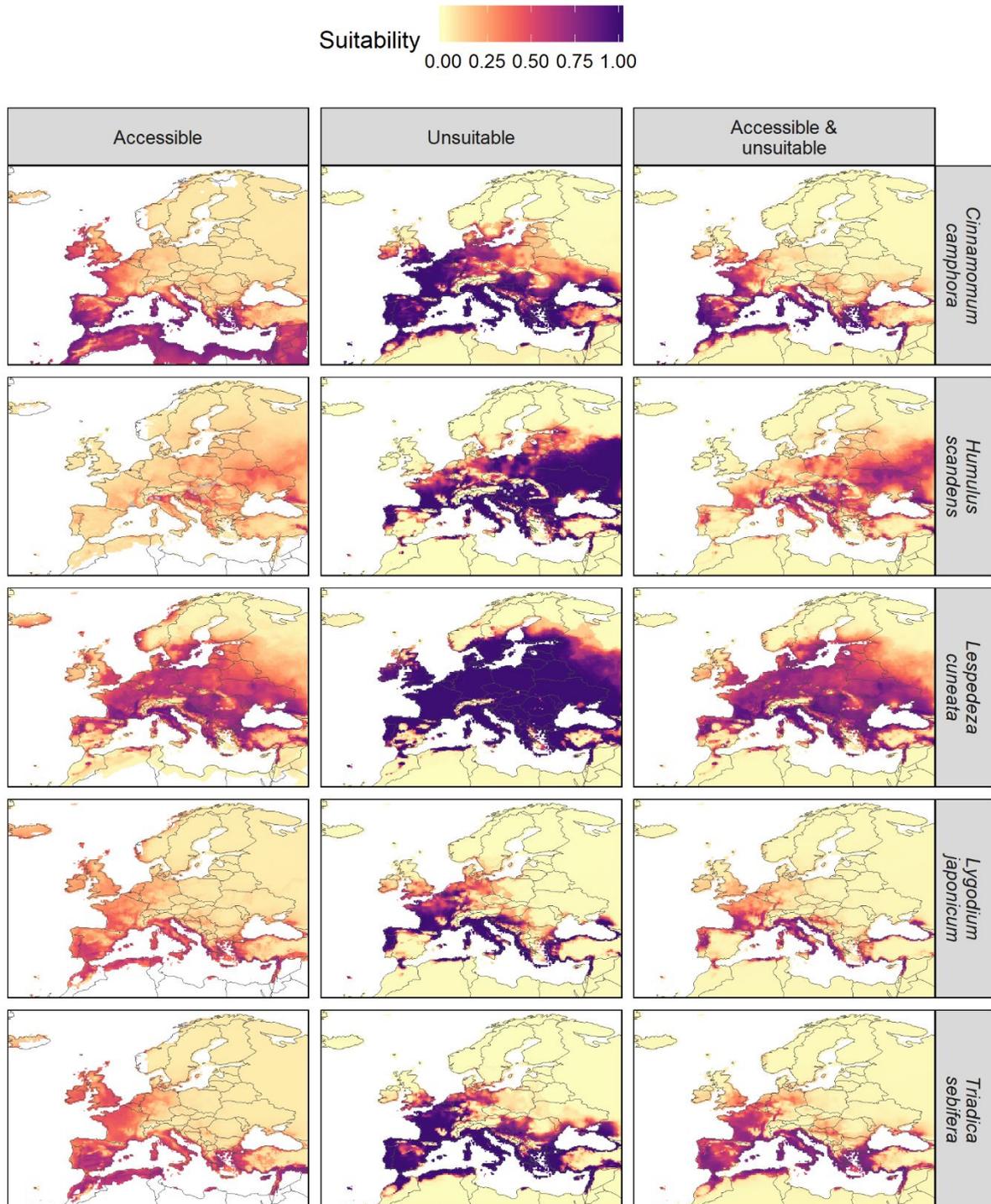


1  
2  
3 429 **Figure 4.** Potential non-native distributions of five Asian plant species in the USA, where all are already  
4 established invasive non-native species with expanding ranges. Projections are from models where the  
5 430 background domain is either just the accessible area, just the unsuitable area or the combined accessible  
6 431 background domain is either just the accessible area, just the unsuitable area or the combined accessible  
7 432 and unsuitable region. Grey points show the species occurrences. Blank land areas are where the model  
8 433 could not project suitability because one or more predictors was outside the range of the training data.



434

435 **Figure 5.** Potential distributions of five Asian plant species in Europe, where the species are currently  
 436 absent or emerging invasive non-native species, equivalent to Figure 4.



437

1  
2  
3 438 **References**  
4  
5

- 6 439 Acevedo, P., Jiménez-Valverde, A., Lobo, J.M., & Real, R. (2012) Delimiting the  
7  
8 440 geographical background in species distribution modelling. *Journal of Biogeography*,  
9  
10 441 **39**, 1383–1390.  
11  
12  
13 442 Austin, M. (2002) Spatial prediction of species distribution: an interface between ecological  
14  
15 443 theory and statistical modelling. *Ecological Modelling*, **157**, 101–118.  
16  
17  
18 444 Barbet-Massin, M., Jiguet, F., Albert, C.H., & Thuiller, W. (2012) Selecting pseudo-absences  
19  
20 445 for species distribution models: how, where and how many? *Methods in Ecology and*  
21  
22 446 *Evolution*, **3**, 327–338.  
23  
24  
25  
26 447 Barrilleaux, T.C. & Grace, James, B. (2000) Growth and invasive potential of *Sapium*  
27  
28 448 *sebiferum* (Euphorbiaceae) within the coastal prairie region: the effects of soil and  
29  
30 449 moisture regime. *American Journal of Botany*, **87**, 1099–1106.  
31  
32  
33  
34 450 Barve, N., Barve, V., Jiménez-Valverde, A., Lira-Noriega, A., Maher, S.P., Peterson, A.T.,  
35  
36 451 Soberón, J., & Villalobos, F. (2011) The crucial role of the accessible area in ecological  
37  
38 452 niche modeling and species distribution modeling. *Ecological Modelling*, **222**, 1810–  
39  
40 453 1819.  
41  
42  
43  
44 454 Boakes, E.H., McGowan, P.J.K., Fuller, R.A., Chang-Qing, D., Clark, N.E., O’Connor, K., &  
45  
46 455 Mace, G.M. (2010) Distorted views of biodiversity: Spatial and temporal bias in species  
47  
48 456 occurrence data. *PLoS Biology*, **8**, e1000385.  
49  
50  
51 457 Bower, M.J., Aslan, C.E., & Rejmánek, M. (2009) Invasion potential of Chinese tallowtree  
52  
53 458 (*Triadica sebifera*) in California’s Central Valley. *Invasive Plant Science and*  
54  
55 459 *Management*, **2**, 386–395.  
56  
57  
58  
59 460 Branquart, E., Brundu, G., Buholzer, S., Chapman, D., Ehret, P., Fried, G., Starfinger, U., van  
60

- 1  
2  
3 461 Valkenburg, J., & Tanner, R. (2016) A prioritization process for invasive alien plant  
4  
5 462 species incorporating the requirements of EU Regulation no. 1143/2014. *EPPO Bulletin*,  
6  
7 463 **46**, 603–617.  
8  
9  
10 464 Broennimann, O. & Guisan, A. (2008) Predicting current and future biological invasions:  
11  
12 both native and invaded ranges matter. *Biology Letters*, **4**, 585–589.  
13  
14  
15 466 Broennimann, O., Treier, U.A., Müller-Schärer, H., Thuiller, W., Peterson, A.T., & Guisan,  
16  
17 A. (2007) Evidence of climatic niche shift during biological invasion. *Ecology Letters*,  
18  
19 467 **10**, 701–709.  
20  
21 468  
22  
23 469 Chapman, D.S., Haynes, T., Beal, S., Essl, F., & Bullock, J.M. (2014) Phenology predicts the  
24  
25 native and invasive range limits of common ragweed. *Global Change Biology*, **20**, 192–  
26  
27 470 202.  
28  
29  
30 472 Chapman, D.S., Makra, L., Albertini, R., Bonini, M., Páldy, A., Rodinkova, V., Šikoparija,  
31  
32 B., Weryszko-Chmielewska, E., & Bullock, J.M. (2016) Modelling the introduction and  
33  
34 473 spread of non-native species: international trade and climate change drive ragweed  
35  
36 474 invasion. *Global change biology*, **22**, 3067–3079.  
37  
38 475  
39  
40 476 Chapman, D.S., Scalone, R., Štefanić, E., & Bullock, J.M. (2017) Mechanistic species  
41  
42 477 distribution modeling reveals a niche shift during invasion. *Ecology*, **98**, 1671–1680.  
43  
44  
45 478 Chefaoui, R.M. & Lobo, J.M. (2007) Assessing the effects of pseudo-absences on predictive  
46  
47 479 distribution model performance. *Ecological Modelling*, **210**, 478–486.  
48  
49  
50 480 Crosby, T. (1993) *How to Detect and Handle Outliers*. ASOC Quality Press, Milwaukee.  
51  
52  
53 481 Dormann, C.F., Elith, J., Bacher, S., Buchmann, C., Carl, G., Carré, G., Marquéz, J.R.G.,  
54  
55 482 Gruber, B., Lafourcade, B., Leitão, P.J., Münkemüller, T., McClean, C., Osborne, P.E.,  
56  
57 483 Reineking, B., Schröder, B., Skidmore, A.K., Zurell, D., & Lautenbach, S. (2012)  
58  
59  
60

- 1  
2  
3 484 Collinearity: a review of methods to deal with it and a simulation study evaluating their  
4  
5 485 performance. *Ecography*, **36**, 27–46.  
6  
7  
8 486 Dormann, C.F., Schymanski, S.J., Cabral, J., Chuine, I., Graham, C., Hartig, F., Kearney, M.,  
9  
10 487 Morin, X., Römermann, C., Schröder, B., & Singer, A. (2011) Correlation and process  
11  
12 488 in species distribution models: bridging a dichotomy. *Journal of Biogeography*, **39**,  
13  
14 489 2119–2131.  
15  
16  
17  
18 490 Elith, J. (2013) Predicting distributions of invasive species. *Invasive Species: Risk Assessment*  
19  
20 491 *and Management* (ed. by A.P. Robinson, T. Walshe, M.A. Burgman, and M. Nunn), pp.  
21  
22 492 93–129. Cambridge University Press, Cambridge, UK.  
23  
24  
25  
26 493 Elith, J., Kearney, M., & Phillips, S. (2010) The art of modelling range-shifting species.  
27  
28 494 *Methods in Ecology and Evolution*, **1**, 330–342.  
29  
30  
31 495 Fitzpatrick, M.C. & Hargrove, W.W. (2009) The projection of species distribution models  
32  
33 496 and the problem of non-analog climate. *Biodiversity and Conservation*, **18**, 2255–2261.  
34  
35  
36 497 Gallien, L., Münkemüller, T., Albert, C.H., Boulangeat, I., & Thuiller, W. (2010) Predicting  
37  
38 498 potential distributions of invasive species: where to go from here? *Diversity and*  
39  
40 499 *Distributions*, **16**, 331–342.  
41  
42  
43  
44 500 Gormley, A.M., Forsyth, D.M., Griffioen, P., Lindeman, M., Ramsey, D.S.L., Scroggie,  
45  
46 501 M.P., & Woodford, L. (2011) Using presence-only and presence-absence data to  
47  
48 502 estimate the current and potential distributions of established invasive species. *Journal*  
49  
50 503 *of Applied Ecology*, **48**, 25–34.  
51  
52  
53  
54 504 Gucker, C. (2010) Available at: <http://www.fs.fed.us/database/feis/>.  
55  
56  
57 505 Guevara, L., Gerstner, B.E., Kass, J.M., & Anderson, R.P. (2017) Toward ecologically  
58  
59 506 realistic predictions of species distributions: A cross-time example from tropical  
60

- 1  
2  
3 507 montane cloud forests. *Global Change Biology*, **24**, 1511–1522.
- 4  
5  
6 508 Guillera-Arroita, G., Lahoz-Monfort, J.J., Elith, J., Gordon, A., Kujala, H., Lentini, P.E.,  
7  
8 509 Mccarthy, M.A., Tingley, R., & Wintle, B.A. (2015) Is my species distribution model fit  
9  
10 510 for purpose? Matching data and models to applications. *Global Ecology and*  
11  
12 511 *Biogeography*, **24**, 276–292.
- 13  
14  
15  
16 512 Hijmans, R.J., Cameron, S.E., Parra, J.L., Jones, P.G., Jarvis, A., Hijmans, R.J., Cameron,  
17  
18 513 S.E., Parra, J.L., Jones, P.G., & Jarvis, A. (2005) Very high resolution interpolated  
19  
20 514 climate surfaces for global land areas, Very high resolution interpolated climate surfaces  
21  
22 515 for global land areas. *International Journal of Climatology*, **25**, 1965–1978.
- 23  
24  
25  
26 516 Jeschke, J.M. & Strayer, D.L. (2008) Usefulness of bioclimatic models for studying climate  
27  
28 517 change and invasive species. *Annals of the New York Academy of Sciences*, **1134**, 1–24.
- 29  
30  
31 518 Jiménez-Valverde, A., Peterson, A.T., Soberón, J., Overton, J.M., Aragón, P., & Lobo, J.M.  
32  
33 519 (2011) Use of niche models in invasive species risk assessments. *Biological Invasions*,  
34  
35 520 **13**, 2785–2797.
- 36  
37  
38  
39 521 Keller, R.P., Lodge, D.M., & Finnoff, D.C. (2007) Risk assessment for invasive species  
40  
41 522 produces net bioeconomic benefits. *Proceedings of the National Academy of Sciences*,  
42  
43 523 **104**, 203–207.
- 44  
45  
46 524 Lobo, J. (2008) AUC : A misleading measure of the performance of predictive distribution  
47  
48 525 models. *Global ecology and Biogeography*, **17**, 145–151.
- 49  
50  
51 526 Mainali, K.P., Warren, D.L., Dhileepan, K., Mcconnachie, A., Strathie, L., Hassan, G., Karki,  
52  
53 527 D., Shrestha, B.B., & Parmesan, C. (2015) Projecting future expansion of invasive  
54  
55 528 species: Comparing and improving methodologies for species distribution modeling.  
56  
57 529 *Global Change Biology*, **21**, 4464–4480.
- 58  
59  
60

- 1  
2  
3 530 Le Maitre, D.C., Thuiller, W., & Schonegevel, L. (2008) Developing an approach to defining  
4  
5 531 the potential distributions of invasive plant species: A case study of *Hakea* species in  
6  
7 532 South Africa. *Global Ecology and Biogeography*, **17**, 569–584.  
8  
9  
10 533 Pearce, J.L. & Boyce, M.S. (2006) Modelling distribution and abundance with presence-only  
11  
12 534 data. *Journal of Applied Ecology*, **43** SRC-, 405–412.  
13  
14  
15 535 Peterson, A.T. & Robins, C.R. (2003) Using ecological-niche modeling to predict barred owl  
16  
17 536 invasions with implications for spotted owl conservation. *Conservation Biology*, **17**,  
18  
19 537 1161–1165.  
20  
21  
22  
23 538 Phillips, S.J. (2009) Sample selection bias and presence-only distribution models:  
24  
25 539 implications for background and pseudo-absence data. *Ecological Applications*, **19**,  
26  
27 540 181–197.  
28  
29  
30  
31 541 Phillips, S.J., Dudik, M., Dudik, M., & Phillips, S.J. (2008) Modeling of species distributions  
32  
33 542 with Maxent: new extensions and a comprehensive evaluation. *Ecography*, **31**, 161–175.  
34  
35  
36 543 Ranc, N., Santini, L., Rondinini, C., Boitani, L., Poitevin, F., Angerbjörn, A., & Maiorano, L.  
37  
38 544 (2017) Performance tradeoffs in target-group bias correction for species distribution  
39  
40 545 models. *Ecography*, **40**, 1076–1087.  
41  
42  
43  
44 546 Roy, H.E., Peyton, J., Aldridge, D.C., et al. (2014) Horizon scanning for invasive alien  
45  
46 547 species with the potential to threaten biodiversity in Great Britain. *Global Change*  
47  
48 548 *Biology*, **20**, 3859–3871.  
49  
50  
51 549 Senay, S.D., Worner, S.P., & Ikeda, T. (2013) Novel three-step pseudo-absence selection  
52  
53 550 technique for improved species distribution modelling. *PLoS One*, **8**, e71218.  
54  
55  
56  
57 551 Siefert, A., Ravenscroft, C., Althoff, D., Alvarez-Yépez, J.C., Carter, B.E., Glennon, K.L.,  
58  
59 552 Heberling, J.M., Jo, I.S., Pontes, A., Sauer, A., Willis, A., & Fridley, J.D. (2012) Scale  
60

- 1  
2  
3 553 dependence of vegetation-environment relationships: A meta-analysis of multivariate  
4  
5 554 data. *Journal of Vegetation Science*, **23**, 942–951.  
6  
7  
8 555 Storkey, J., Stratonovitch, P., Chapman, D.S., Vidotto, F., & Semenov, M.A. (2014) A  
9  
10 556 process-based approach to predicting the effect of climate change on the distribution of  
11  
12 557 an invasive allergenic plant in europe. *Plos One*, **9**, e88156.  
13  
14  
15 558 Talluto, M. V, Boulangeat, I., Ameztegui, A., Aubin, I., Berteaux, D., Butler, A., Doyon, F.,  
16  
17 559 Drever, C.R., Fortin, M.-J., Franceschini, T., Liénard, J., McKenney, D., Solarik, K.A.,  
18  
19 560 Strigul, N., Thuiller, W., & Gravel, D. (2015) Cross-scale integration of knowledge for  
20  
21 561 predicting species ranges: a metamodeling framework. *Global Ecology and*  
22  
23 562 *Biogeography*, **25**, 238–249.  
24  
25  
26  
27  
28 563 Tanner, R., Branquart, E., Brundu, G., Buholzer, S., Chapman, D., Ehret, P., Fried, G.,  
29  
30 564 Starfinger, U., & van Valkenburg, J. (2017) The prioritisation of a short list of alien  
31  
32 565 plants for risk analysis within the framework of the Regulation (EU) No. 1143/2014.  
33  
34 566 *NeoBiota*, **35**, 87–118.  
35  
36  
37  
38 567 Thuiller, W., Brotons, L., Araújo, M.B., & Lavorel, S. (2004) Effects of restricting  
39  
40 568 environmental range of data to project current and future species distributions.  
41  
42 569 *Ecography*, **27**, 165–172.  
43  
44  
45 570 Thuiller, W., Georges, D., Engler, R., & Breiner, F. (2016) biomod2: Ensemble platform for  
46  
47 571 species distribution modeling. R package version 3.3-7. Available at: [https://cran.r-](https://cran.r-project.org/web/packages/biomod2/index.html)  
48  
49 572 [project.org/web/packages/biomod2/index.html](https://cran.r-project.org/web/packages/biomod2/index.html), .  
50  
51  
52  
53 573 Thuiller, W., Lafourcade, B., Engler, R., & Araújo, M.B. (2009) BIOMOD - A platform for  
54  
55 574 ensemble forecasting of species distributions. *Ecography*, **32**, 369–373.  
56  
57  
58 575 Václavík, T. & Meentemeyer, R.K. (2009) Invasive species distribution modeling (iSDM):  
59  
60

- 1  
2  
3 576 Are absence data and dispersal constraints needed to predict actual distributions?  
4  
5 577 *Ecological Modelling*, **220**, 3248–3258.  
6  
7  
8 578 VanDerWal, J., Shoo, L.P., Graham, C., & Williams, S.E. (2009) Selecting pseudo-absence  
9  
10 579 data for presence-only distribution modeling: How far should you stray from what you  
11  
12 580 know? *Ecological Modelling*, **220**, 589–594.  
13  
14  
15  
16 581 Vilà, M., Espinar, J.L., Hejda, M., Hulme, P.E., Jarošík, V., Maron, J.L., Pergl, J., Schaffner,  
17  
18 582 U., Sun, Y., & Pyšek, P. (2011) Ecological impacts of invasive alien plants: A meta-  
19  
20 583 analysis of their effects on species, communities and ecosystems. *Ecology Letters*, **14**,  
21  
22 584 702–708.  
23  
24  
25  
26 585 Warton, D.I. & Shepherd, L.C. (2010) Poisson point process models solve the “pseudo-  
27  
28 586 absence problem” for presence-only data in ecology. *Annals of Applied Statistics*, **4**,  
29  
30 587 1383–1402.  
31  
32  
33  
34 588 Zhao, X., Wang, G., Shen, Z., Zhang, H., & Qiu, M. (2006) Impact of elevated CO<sub>2</sub>  
35  
36 589 concentration under three soil water levels on growth of *Cinnamomum camphora*.  
37  
38 590 *Journal of Zhejiang University, Science B*, **7**, 283–290.  
39  
40  
41 591 Zomer, R.J., Trabucco, A., Bossio, D.A., & Verchot, L. V (2008) Climate change mitigation:  
42  
43 592 A spatial analysis of global land suitability for clean development mechanism  
44  
45 593 afforestation and reforestation. *Agr Ecosyst Environ*, **126**, 67–80.  
46  
47  
48  
49 594  
50  
51 595  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 596 **Data Accessibility**  
4  
5

6 597 A data file containing the 0.25 x 0.25 gridded data on climate, recording effort, species  
7  
8 598 occurrence, accessibility and unsuitability is included in the Supporting Information.  
9

10  
11 599 **Biosketch**  
12

13  
14 600 The research team focuses on risk assessment for emerging invasive non-native species in Europe.  
15

16 601 Among other factors contributing to risk, the team use global-scale species distribution modelling to  
17

18 602 identify the suitable conditions for establishment by the focal species and use this to project their  
19

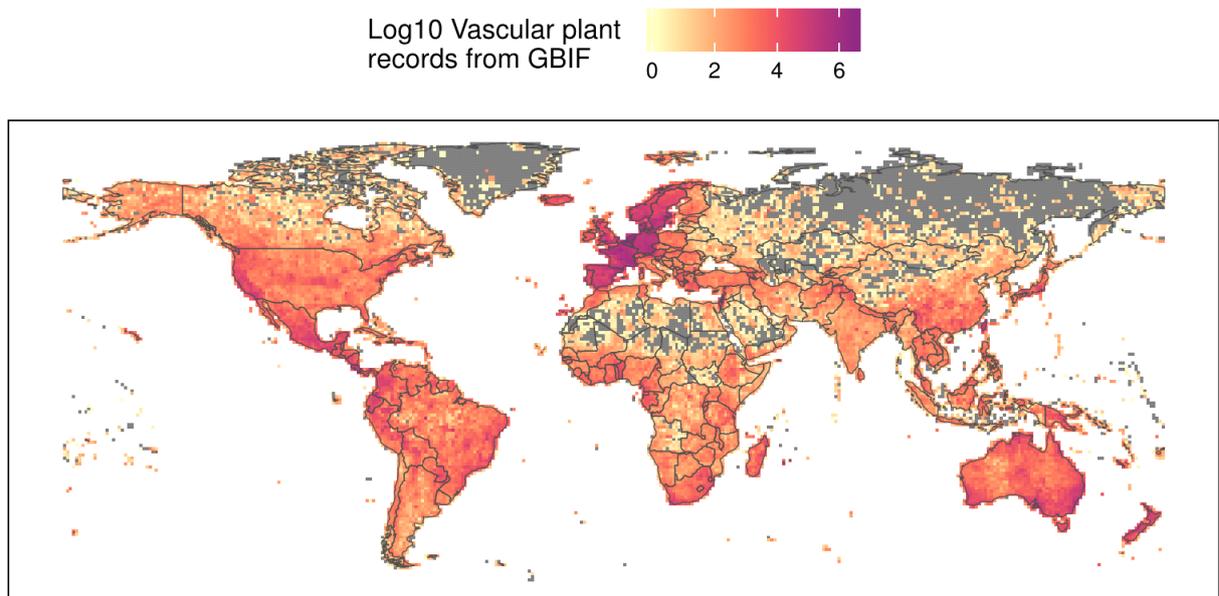
20 603 potential distributional range in the risk assessment area.  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

**Supporting Information: Improving species distribution models for invasive non-native species with biologically informed pseudo-absence selection**

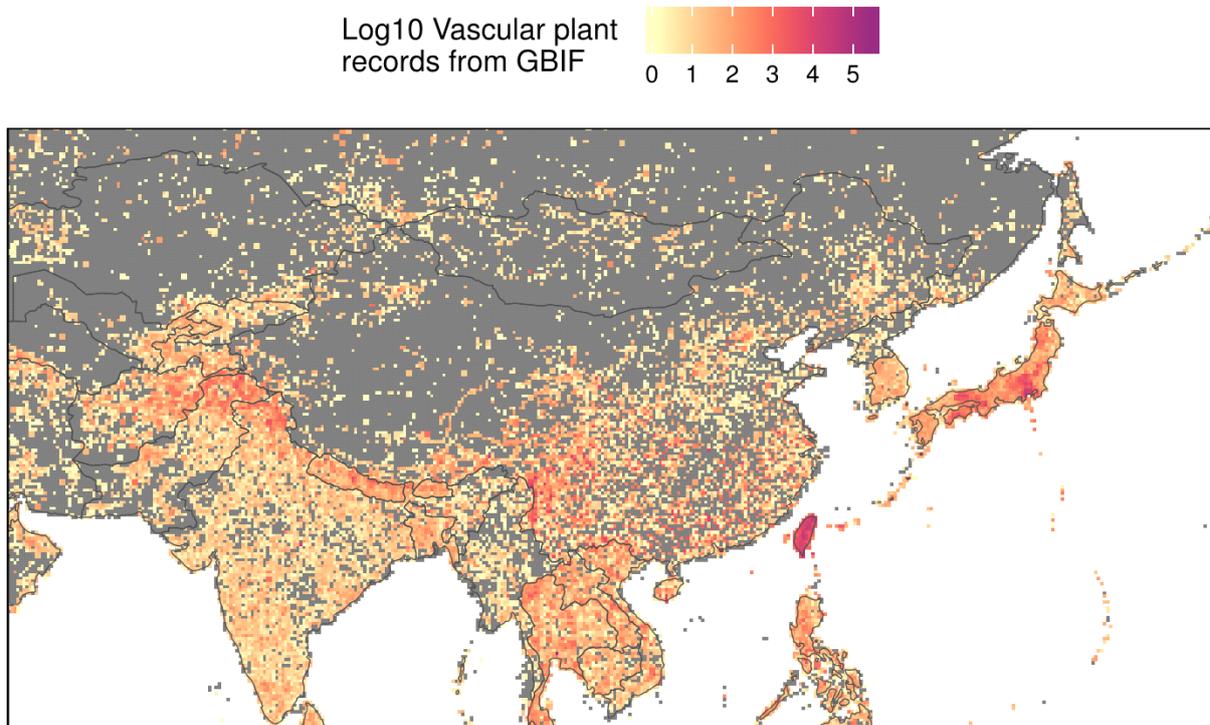
Daniel Chapman, Oliver L. Pescott, Helen E. Roy, Rob Tanner

**Appendix S1 – Proxy for recording effort**

**Figure S1.1.** The global density of vascular plant (phylum Tracheophyta) records retrieved from the Global Biodiversity Information Facility, mapped on a 1x1 degree grid and displayed on a log<sub>10</sub> scale. Dark grey areas returned no records.

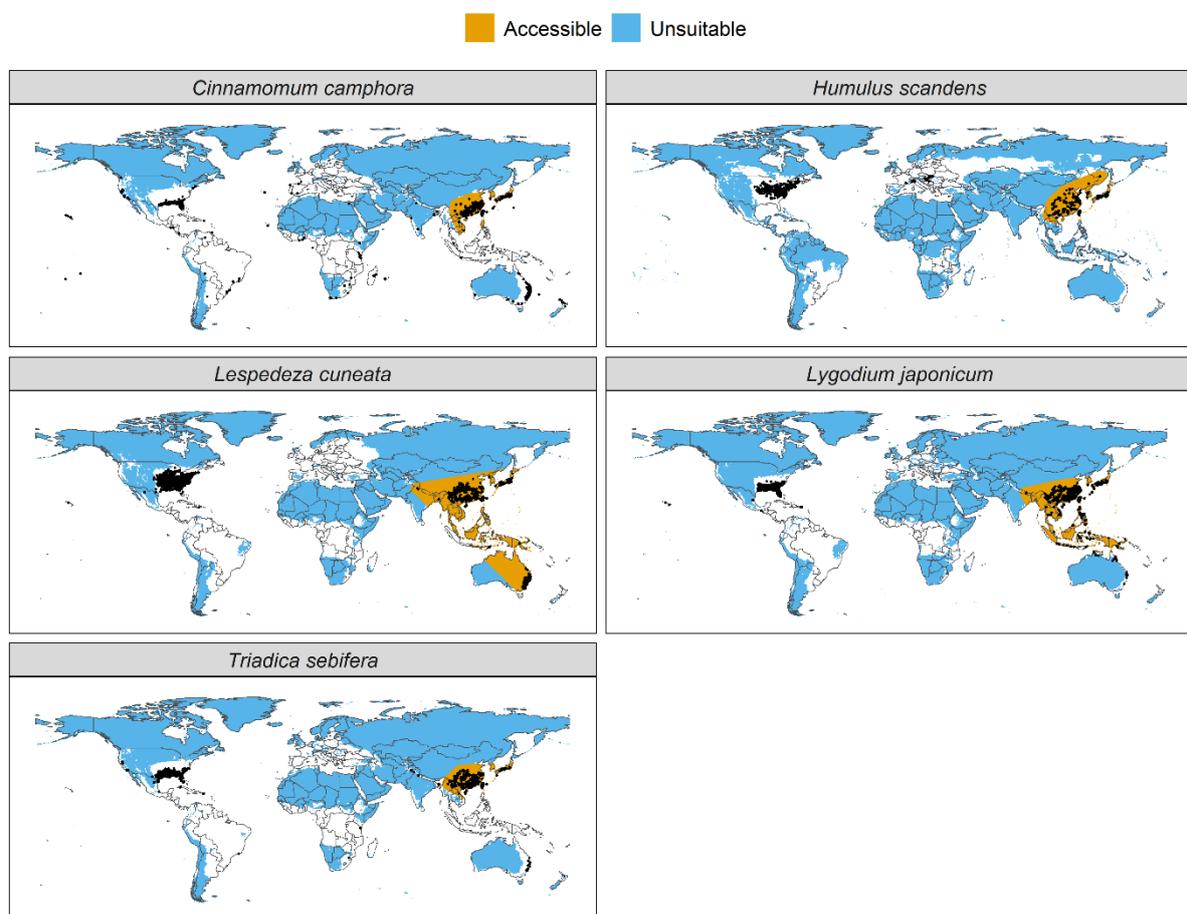


**Figure S1.2.** The density of vascular plant (phylum Tracheophyta) records retrieved from the Global Biodiversity Information Facility at the spatial resolution of the model (0.25 x 0.25 degree grid) for Asia. This region spanned the native range of most of the study species and had pronounced recording bias. Dark grey areas returned no records.



## Appendix S2 – Distributions of the five modelled plant species and definitions of environmentally-unsuitable regions

**Figure S2.3.** Distribution records of the five species that were modelled, plotted as black points. Background shading shows the accessible and unsuitable background domains from which ‘pseudo-absences’ were drawn. The accessible background is only visible in the native range of the species, as it is otherwise masked by the non-native presences. Our approach excludes parts of the world that may be environmentally suitable, but are out of dispersal range (unshaded).



**Table S2.1.** Rules for defining highly unsuitable conditions for establishment by the five study species, with respect to three climatic variables. The rules are based on a combination of prior biological knowledge about key constraints and the extreme values of the climate variables at the species occurrence locations. For CMI and upper limit on Bio10, few direct estimates were used, so the 0.5<sup>th</sup> or 99.5<sup>th</sup> percentile values at occurrence locations were used as thresholds, respectively. The individual rules were combined with an OR statement to generate unsuitable background domains for the distribution modelling.

<b>Species</b>	<b>Bio6 (mean minimum temperature during coldest month)</b>	<b>Bio10 (mean temperature during warmest quarter)</b>	<b>CMI (climatic moisture index)</b>
<i>Cinnamomum camphora</i>	<-10 °C; the temperature causing frost damage in overwintering seedlings (You et al., 2008).	<15 °C; the reported minimum annual temperature (Orwa et al., 2009; CABI, 2018), <b>OR</b> >29 °C	<0.25
<i>Humulus scandens</i>	>16 °C; overwintering seeds require stratification (Balogh & Dancza, 2008) and the warmest occurrence is at 15.5 °C.	<15 °C; approximately corresponds to a known requirement of ~1300 degree days (base 4 °C) for maturation (G. Fried, personal communication), <b>OR</b> >28 °C	<0.45
<i>Lespedeza cuneata</i>	<-17 °C; effects of frost are uncertain but may contribute to mortality (Gucker, 2010) and the coldest occurrence is at 16.4 °C.	<13 °C; consistent with minimum temperatures for germination (Qiu et al., 1995) and seedling growth (Hill & Luck, 1991) , <b>OR</b> >29 °C; consistent with observed reductions in leaf size and height when grown above 25 °C (Kalburtji et al., 2007)	<0.45
<i>Lygodium japonicum</i>	<-8 °C; consistent with mean values of Bio6 in USDA Plant Hardiness Zone 6, where the species is semi-hardy (Loan, 2006).	<16 °C; no information available on summer temperature requirements, so the 0.5 <sup>th</sup> percentile value of occurrences was assumed to be a low limit, <b>OR</b> >29 °C	<0.55
<i>Triadica sebifera</i>	<-9 °C; frost is considered the strongest factor limiting invasion in the USA (Gan et al., 2009) and the coldest occurrence is at -8.5 °C.	< 16 °C; consistent with the temperature inhibiting seed germination (Nijjer et al., 2002) , <b>OR</b> >29 °C	<0.30

## Appendix S3 – Model summary tables

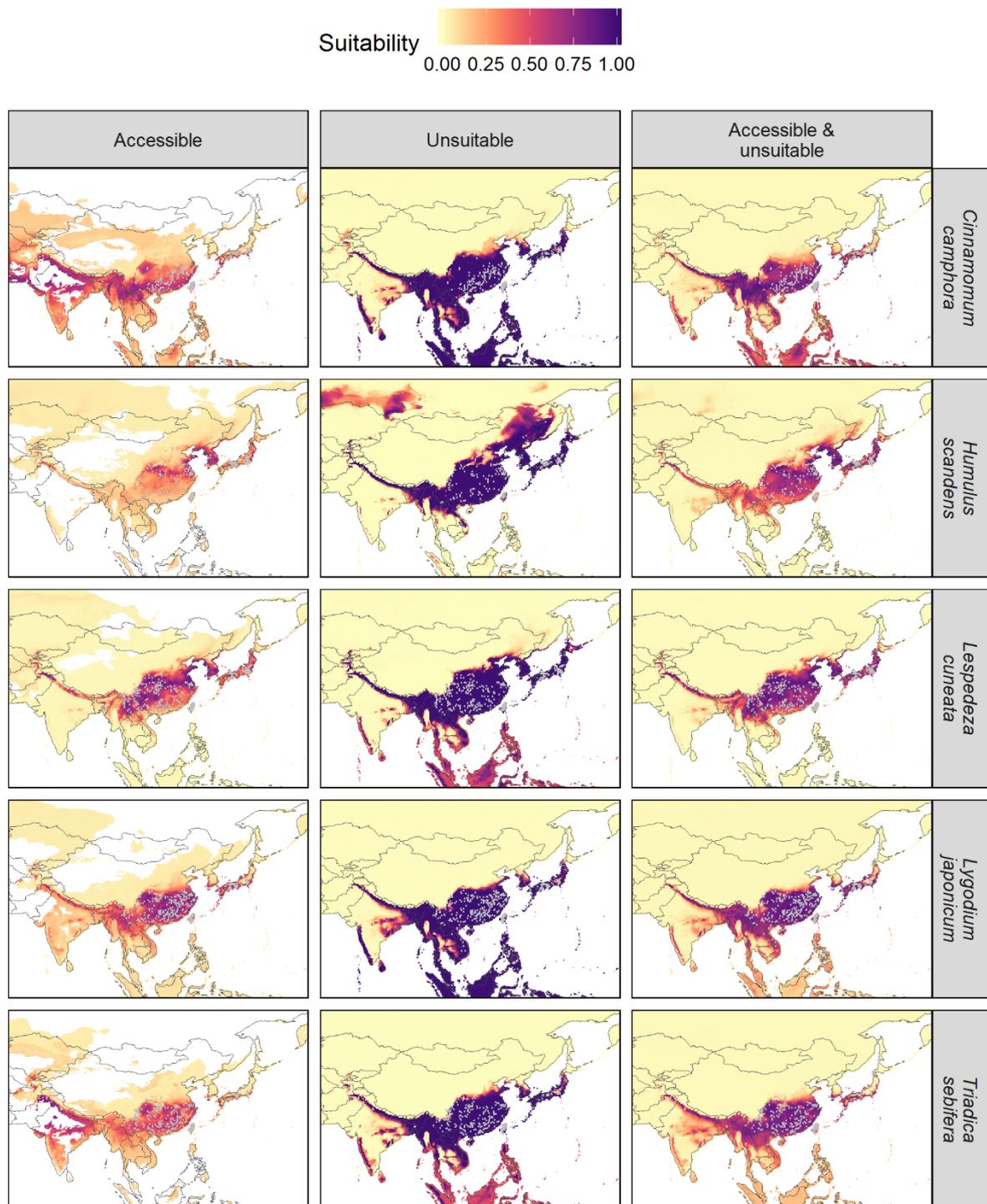
**Table S3.2.** Details of the ensemble presence-background models for five Asian plant species occurring elsewhere as invasive non-native species. Three model specifications were applied, where the background sample was drawn from only an accessible background, only from an unsuitable background or from a combined accessible and unsuitable background. Seven individual algorithms were fitted and evaluated based on their AUC, cross-validated within the training data. Poorly performing algorithms were rejected, and the remaining ones combined into the ensemble. Variable importances are given as percentages for Bio6 (mean minimum temperature of coldest month), Bio10 (mean temperature of warmest quarter) and CMI (ratio of annual precipitation to potential evapotranspiration).

Species	Algorithm	Accessible background					Unsuitable background					Accessible and unsuitable background				
		AUC	In ensemble	Bio6	Bio10	CMI	AUC	In ensemble	Bio6	Bio10	CMI	AUC	In ensemble	Bio6	Bio10	CMI
<i>Cinnamomum camphora</i>	ANN	0.758	yes	74%	12%	13%	0.998	yes	39%	31%	30%	0.946	yes	50%	24%	26%
	GAM	0.753	yes	93%	4%	3%	0.997	yes	43%	29%	28%	0.941	no	45%	30%	24%
	GBM	0.771	yes	83%	8%	8%	0.992	no	42%	33%	25%	0.952	yes	48%	32%	21%
	GLM	0.750	yes	89%	6%	4%	0.996	yes	43%	29%	28%	0.937	no	51%	26%	24%
	MARS	0.750	yes	91%	3%	6%	0.996	yes	44%	28%	28%	0.948	yes	49%	28%	23%
	Maxent	0.749	yes	77%	9%	14%	0.998	yes	28%	44%	27%	0.952	yes	49%	31%	21%
	RF	0.726	no	52%	23%	24%	0.999	yes	36%	34%	30%	0.948	yes	48%	27%	25%
	<b>Ensemble</b>	<b>0.769</b>		<b>85%</b>	<b>7%</b>	<b>8%</b>	<b>0.999</b>		<b>39%</b>	<b>33%</b>	<b>29%</b>	<b>0.954</b>		<b>49%</b>	<b>28%</b>	<b>23%</b>
<i>Humulus scandens</i>	ANN	0.773	yes	55%	20%	25%	0.999	yes	24%	41%	34%	0.955	yes	34%	33%	33%
	GAM	0.773	yes	57%	19%	24%	0.998	yes	23%	44%	33%	0.957	yes	34%	33%	33%
	GBM	0.785	yes	54%	16%	30%	0.999	yes	22%	48%	30%	0.958	yes	37%	27%	36%
	GLM	0.765	yes	46%	26%	28%	0.999	yes	27%	40%	33%	0.953	yes	32%	35%	33%
	MARS	0.771	yes	56%	18%	27%	1.000	yes	22%	45%	33%	0.956	yes	34%	31%	35%
	Maxent	0.762	yes	54%	20%	26%	0.995	no	18%	48%	34%	0.955	yes	40%	27%	32%
	RF	0.762	yes	42%	26%	32%	0.998	yes	22%	44%	34%	0.953	yes	38%	28%	34%
	<b>Ensemble</b>	<b>0.784</b>		<b>52%</b>	<b>21%</b>	<b>27%</b>	<b>1.000</b>		<b>23%</b>	<b>44%</b>	<b>33%</b>	<b>0.959</b>		<b>36%</b>	<b>30%</b>	<b>34%</b>

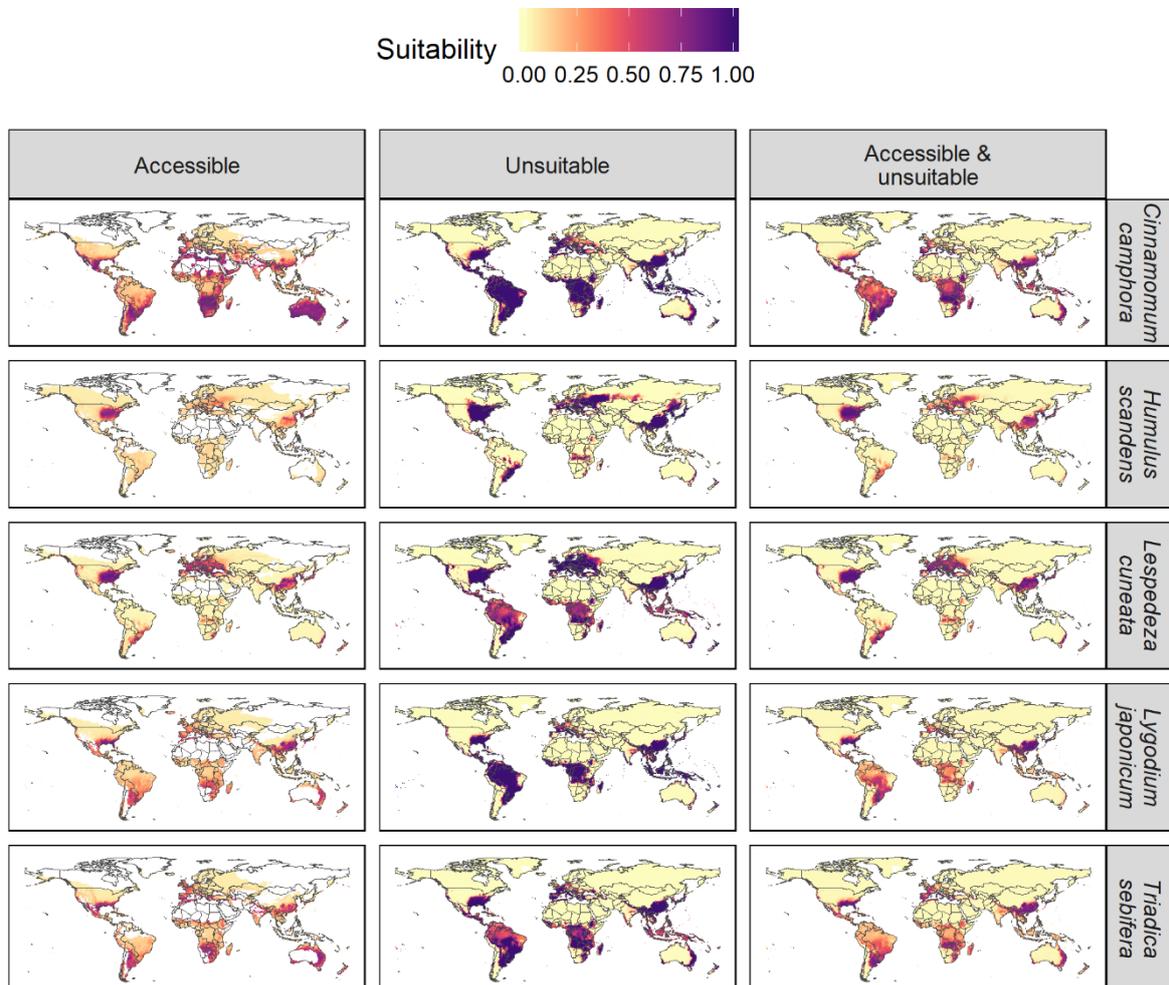
Species	Algorithm	Accessible background					Unsuitable background					Accessible and unsuitable background				
		AUC	In ensemble	Bio6	Bio10	CMI	AUC	In ensemble	Bio6	Bio10	CMI	AUC	In ensemble	Bio6	Bio10	CMI
<i>Lespedeza cuneata</i>	ANN	0.912	yes	49%	11%	40%	0.999	yes	33%	30%	37%	0.963	yes	40%	21%	39%
	GAM	0.910	yes	47%	10%	43%	0.999	yes	37%	24%	39%	0.961	yes	41%	20%	38%
	GBM	0.914	yes	59%	6%	34%	0.999	yes	37%	30%	34%	0.964	yes	43%	19%	38%
	GLM	0.893	no	53%	11%	36%	1.000	yes	36%	25%	39%	0.953	no	42%	20%	38%
	MARS	0.907	yes	48%	11%	41%	0.998	yes	36%	25%	39%	0.963	yes	39%	20%	41%
	Maxent	0.904	yes	40%	18%	42%	0.998	yes	18%	44%	38%	0.962	yes	39%	25%	36%
	RF	0.904	yes	41%	20%	39%	0.999	yes	24%	35%	41%	0.961	yes	36%	25%	38%
	<b>Ensemble</b>	<b>0.912</b>		<b>47%</b>	<b>13%</b>	<b>40%</b>	<b>1.000</b>		<b>32%</b>	<b>30%</b>	<b>38%</b>	<b>0.964</b>		<b>40%</b>	<b>22%</b>	<b>38%</b>
<i>Lygodium japonicum</i>	ANN	0.853	yes	75%	16%	9%	0.998	yes	27%	31%	42%	0.959	yes	41%	26%	33%
	GAM	0.844	yes	89%	9%	2%	0.999	yes	33%	24%	43%	0.958	yes	38%	23%	39%
	GBM	0.857	yes	74%	20%	5%	0.997	yes	34%	31%	35%	0.965	yes	39%	29%	32%
	GLM	0.844	yes	84%	15%	1%	0.999	yes	30%	26%	45%	0.949	no	36%	24%	40%
	MARS	0.850	yes	78%	19%	2%	0.997	yes	28%	28%	44%	0.964	yes	39%	22%	39%
	Maxent	0.841	yes	64%	24%	12%	0.995	no	10%	45%	45%	0.962	yes	39%	30%	31%
	RF	0.838	yes	55%	24%	21%	0.999	yes	20%	36%	44%	0.961	yes	34%	32%	34%
	<b>Ensemble</b>	<b>0.855</b>		<b>74%</b>	<b>18%</b>	<b>8%</b>	<b>0.999</b>		<b>29%</b>	<b>29%</b>	<b>42%</b>	<b>0.966</b>		<b>38%</b>	<b>27%</b>	<b>35%</b>
<i>Triadica sebifera</i>	ANN	0.754	yes	64%	16%	21%	0.999	yes	36%	30%	34%	0.940	yes	44%	24%	32%
	GAM	0.756	yes	78%	13%	9%	1.000	yes	42%	26%	32%	0.933	yes	44%	24%	32%
	GBM	0.768	yes	64%	17%	19%	0.998	yes	40%	30%	30%	0.940	yes	46%	29%	25%
	GLM	0.750	yes	78%	15%	7%	1.000	yes	37%	29%	34%	0.929	no	44%	23%	32%
	MARS	0.750	yes	76%	14%	10%	0.998	yes	39%	29%	31%	0.937	yes	50%	23%	27%
	Maxent	0.755	yes	69%	16%	15%	0.991	no	16%	50%	34%	0.937	yes	45%	31%	24%
	RF	0.719	no	45%	28%	27%	0.998	yes	26%	40%	35%	0.931	no	42%	28%	29%
	<b>Ensemble</b>	<b>0.766</b>		<b>71%</b>	<b>15%</b>	<b>14%</b>	<b>1.000</b>		<b>37%</b>	<b>31%</b>	<b>33%</b>	<b>0.941</b>		<b>46%</b>	<b>26%</b>	<b>28%</b>

## Appendix S4 – Native range and global projections for five modelled plant species

**Figure S4.4.** Projections of suitability for the five study species in Asia, which includes their main native distributions. Plots show outputs from the three background specifications, equivalent to Figure 4 in the main text.



**Figure S4.5.** Global projections of suitability for the five study species. Plots show outputs from the three model specifications for the background region, equivalent to Figure 4 in the main text.



## **Appendix S5 – Sensitivity analysis on model settings**

### *Sensitivity to size of the accessible region and number of background samples*

All five species were modelled as in the main text but with all combinations of:

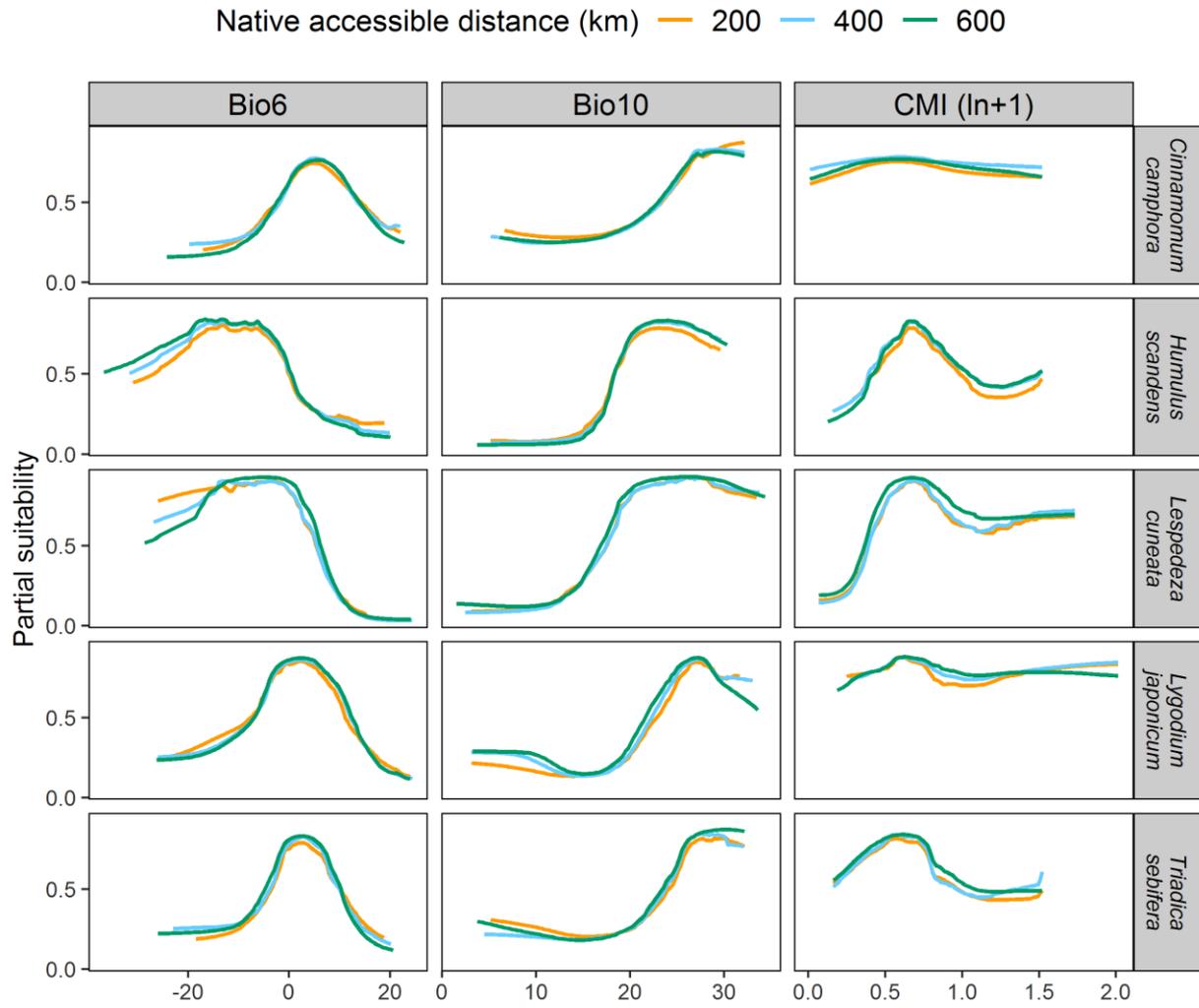
- the buffer radius for the native accessible region set to 0 km (unsuitable area model), 200 km, 400 km (as used in the main text) and 600 km; and
- the number of background samples (pseudo-absences) taken from the unsuitable region set to 0 (accessible area model), 1000, 3000 (as used in the main text) and 5000.

Model response plots were generally not very sensitive to the choice of these settings, except when only 1000 unsuitable background samples were taken (Figures S5.6-5.8). As a result, global suitability projections were almost identical (not shown).

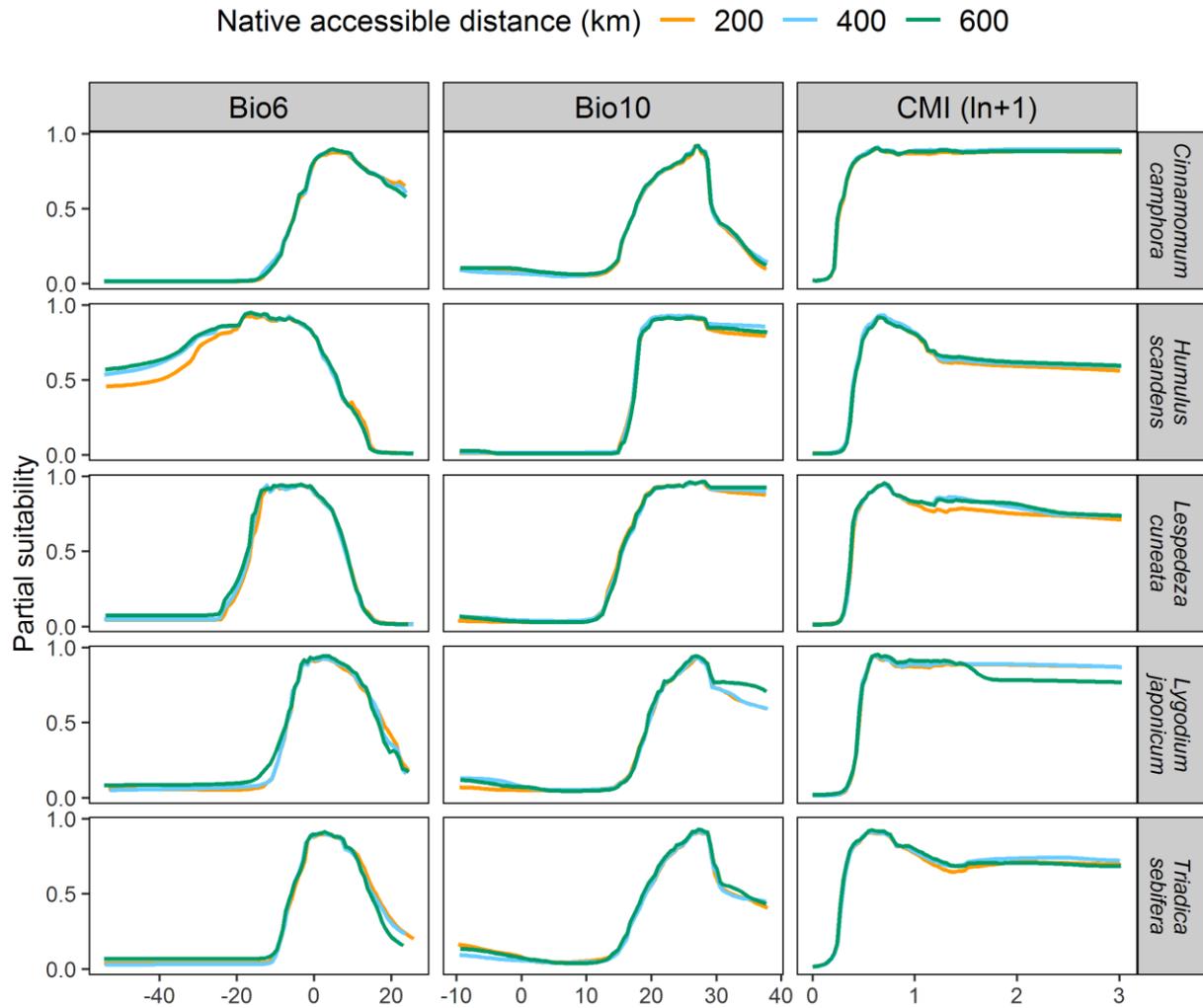
### *Sensitivity to the rules used to define the unsuitable region*

All five species were modelled as in the main text but with the rules defining their unsuitable regions set to either the values in Table S2.1 (used in the main text) or to a more conservative definition of the unsuitable region. For the latter, temperature thresholds were made more extreme by 2 °C (e.g. <-10 °C changed to <-12 °C; >16 °C changed to >18 °C, etc.) and moisture (CMI) thresholds were made more extreme by 10% (e.g. 0.25 changed to 0.15, etc.). The effect of this was to reduce the size of the unsuitable region and separate it more strongly from the species' occurrences. Surprisingly, this generally had little influence on the fitted response functions (Figure S5.9) or projections made from the models (not shown).

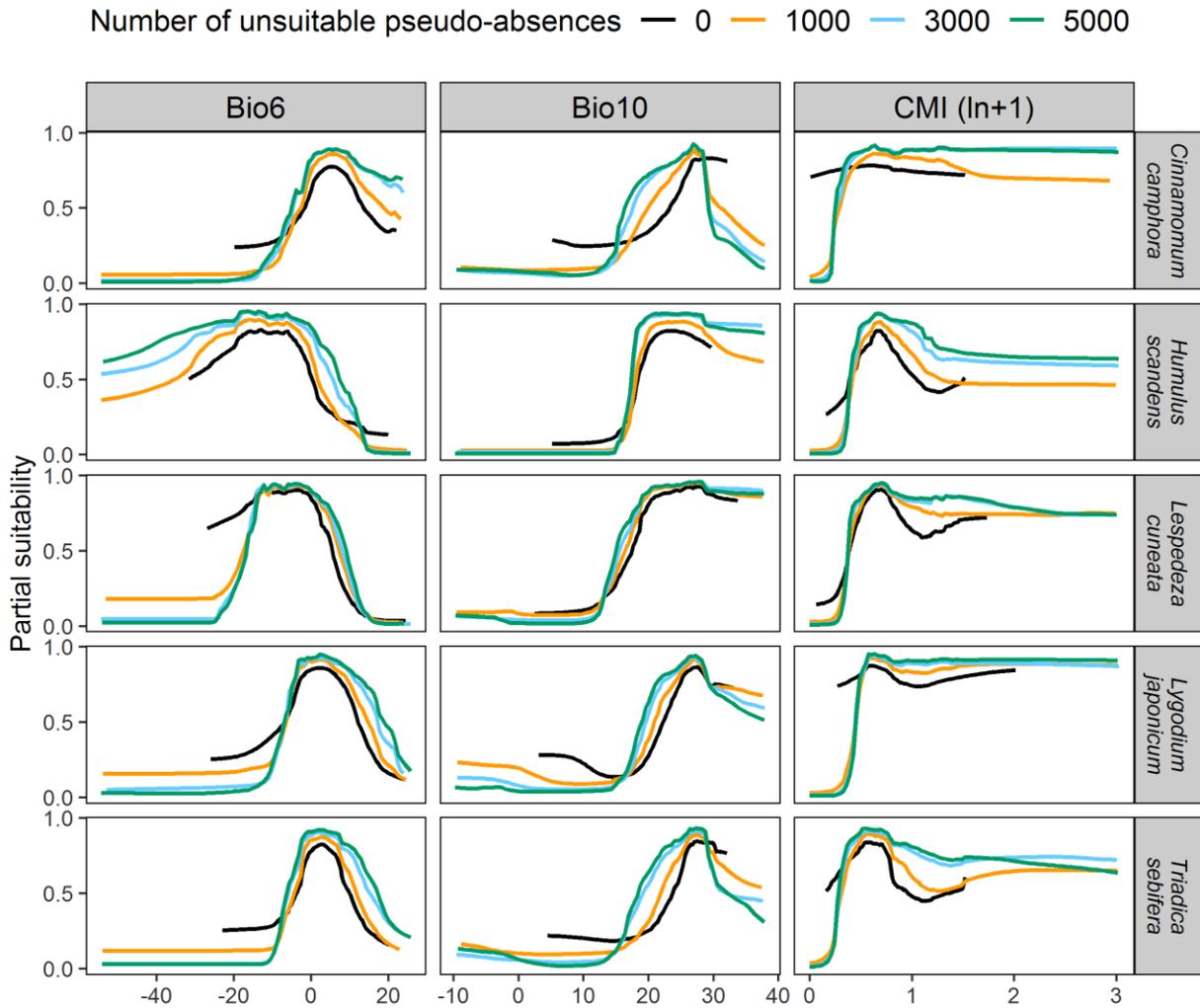
**Figure S5.6.** Partial response plots for models fitted using only accessible backgrounds, and with the native accessible region defined with buffer radii of 200 km, 400 km and 600 km.



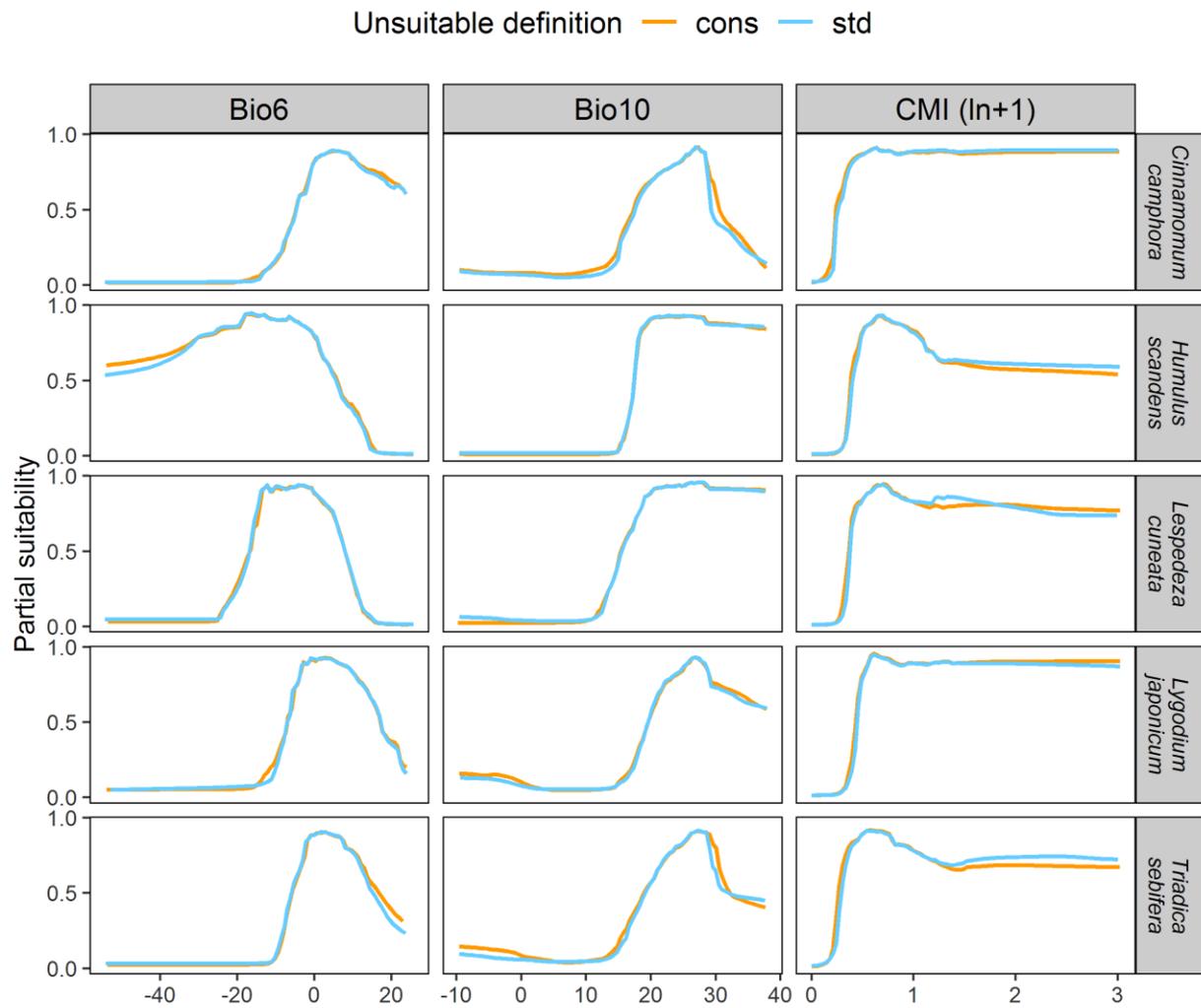
**Figure S5.7.** Partial response plots for models fitted using accessible and unsuitable backgrounds, and with the native accessible region defined with buffer radii of 200 km, 400 km and 600 km and 3000 background samples.



**Figure S5.8.** Effect of the number of unsuitable background samples on response plots fitted by models using accessible and unsuitable backgrounds. Models were fitted with 1000, 3000 or 5000 unsuitable background samples and a 400 km buffer for the native accessible region.



**Figure S5.9.** Response curves fitted to models where the unsuitable region was defined as in the main text (std) or with more conservative rules (cons).



## Appendix S6 – ‘Read me’ for the data file

The file 'INNS SDM data for modelling.rds' contains the data used in the study, compiled from publicly available sources and gridded at 0.25 x 0.25 degree resolution. To open the file in R and read in the data as a data.frame please use the command:

```
readRDS('INNS SDM data for modelling.rds')
```

The .rds file contains a compressed R data.frame with the following columns:

- x = longitude of 0.25 degree grid cell centre
- y = latitude of 0.25 degree grid cell centre
- bio6 = Worldclim Bio6 (Minimum temperature of the coldest quarter, C)
- bio10 = Worldclim Bio10 (Mean temperature of the warmest quarter, C)
- moisture = climatic moisture index (ratio of annual precipitation to potential evapotranspiration)
- effort = proxy for recording effort, the number of Tracheophyte records held by GBIF
- occ.Cinnamomum.camphora = occurrences of *Cinnamomum camphora* (1 = occurrence)
- occ.Humulus.scandens = occurrences of *Humulus scandens* (1 = occurrence)
- occ.Lespedeza.cuneata = occurrences of *Lespedeza cuneata* (1 = occurrence)
- occ.Lygodium.japonicum = occurrences of *Lygodium japonicum* (1 = occurrence)
- occ.Triadica.sebifera = occurrences of *Triadica sebifera* (1 = occurrence)
- native.occ.Cinnamomum.camphora = occurrences of *Cinnamomum camphora* in native range (1 = occurrence)
- native.occ.Humulus.scandens = occurrences of *Humulus scandens* in native range (1 = occurrence)
- native.occ.Lespedeza.cuneata = occurrences of *Lespedeza cuneata* in native range (1 = occurrence)
- native.occ.Lygodium.japonicum = occurrences of *Lygodium japonicum* in native range (1 = occurrence)
- native.occ.Triadica.sebifera = occurrences of *Triadica sebifera* in native range (1 = occurrence)

- nonnative.occ.Cinnamomum.camphora = occurrences of *Cinnamomum camphora* in non-native range (1 = occurrence)
- nonnative.occ.Humulus.scandens = occurrences of *Humulus scandens* in non-native range (1 = occurrence)
- nonnative.occ.Lespedeza.cuneata = occurrences of *Lespedeza cuneata* in non-native range (1 = occurrence)
- nonnative.occ.Lygodium.japonicum = occurrences of *Lygodium japonicum* in non-native range (1 = occurrence)
- nonnative.occ.Triadica.sebifera = occurrences of *Triadica sebifera* in non-native range (1 = occurrence)
- accessible.Cinnamomum.camphora = the accessible background region for *Cinnamomum camphora* (1 = accessible)
- accessible.Humulus.scandens = the accessible background region for *Humulus scandens* (1 = accessible)
- accessible.Lespedeza.cuneata = the accessible background region for *Lespedeza cuneata* (1 = accessible)
- accessible.Lygodium.japonicum = the accessible background region for *Lygodium japonicum* (1 = accessible)
- accessible.Triadica.sebifera = the accessible background region for *Triadica sebifera* (1 = accessible)
- unsuitable.Cinnamomum.camphora = the unsuitable background region for *Cinnamomum camphora* (1 = unsuitable)
- unsuitable.Humulus.scandens = the unsuitable background region for *Humulus scandens* (1 = unsuitable)
- unsuitable.Lespedeza.cuneata = the unsuitable background region for *Lespedeza cuneata* (1 = unsuitable)

- unsuitable.Lygodium.japonicum = the unsuitable background region for *Lygodium japonicum* (1 = unsuitable)
- unsuitable.Triadica.sebifera = the unsuitable background region for *Triadica sebifera* (1 = unsuitable)

## Supporting References

- Balogh, L. & Dancza, I. (2008) *Humulus japonicus*, an emerging invader in Hungary. *Plant invasions: Human perception, ecological impacts and management* (ed. by B. Tokarska-Guzi, J.H. Brock, G. Brundu, C.C. Child, C. Daehler, and P. Pyšek), pp. 73–91. Backhuys Publishers, Leiden, Netherlands.
- CABI (2018) *Invasive species compendium*. CAB International, Wallingford, UK.
- Gan, J., Miller, J.H., Wang, H., & Taylor, J.W. (2009) Invasion of tallow tree into southern US forests: influencing factors and implications for mitigation. *Canadian Journal of Forest Research*, **39**, 1346–1356.
- Gucker, C. (2010) Fire Effects Information System (FEIS). Available at: <https://www.feis-crs.org/feis/>.
- Hill, M.J. & Luck, R. (1991) The effect of temperature on germination and seedling growth of temperate perennial pasture legumes. *Australian Journal of Agricultural Research*, **42**, 175–189.
- Kalbertji, K.L., Mosjidis, J.A., & Mamolos, A.P. (2007) Effects of day-night temperature combinations under constant day length on emergence and early growth of *Sericea lespedeza* genotypes. *Canadian Journal of Plant Science*, **87**, 77–81.
- Loan, A. Van (2006) Japanese climbing fern: The insidious “other” *Lygodium*. *Wildland Weeds*, **9**, 25–27.
- Nijjer, S., Lankau, R.A., Rogers, W.E., & Siemann, E. (2002) Effects of temperature and light on Chinese tallow (*Sapium sebiferum*) and Texas sugarberry (*Celtis laevigata*) seed germination. *Texas Journal of Science*, **54**, 63–68.
- Orwa, C., Mutua, A., Kindt, R., Jamnadass, R., & Anthony S (2009) Tree Functional and Ecological Databases. Available at: <http://www.worldagroforestry.org/output/tree-functional-and-ecological-databases>

Qiu, J., Mosjidis, J.A., & Williams, J.C. (1995) Variability for temperature of germination in *Sericea lespedeza* germplasm. *Crop science*, **35**, 237–241.

You, Y., Liu, H., Wu, R., & Lin, Y. (2008) Effect of low temperature stress on cold resistance of *Cinnamomum camphora* seedling. *Guangdong Agricultural Sciences*, **11**, 23–25.