# Why to choose Random Forest to predict rare species distribution with few samples in large undersampled areas? Three Asian crane species models provide supporting evidence

Chunrong Mi [1,2] , Falk Huettmann [3] , Yumin Guo [Corresp., 1] , Xuesong Han [1] , Lijia Wen [1]

[1] College of Nature Conservation, Beijing Forestry University, Beijing, China

[2] Key Laboratory of Water Cycle and Related Land Surface Processes, Institute of Geographic Sciences and Natural Resources Research, University of Chinese Academy of Sciences, Beijing, China

[3] EWHALE Lab, Department of Biology and Wildlife, Institute of Arctic Biology, University of Alaska Fairbanks (UAF), Fairbanks, Alaska, United States

Corresponding Author: Yumin Guo
Email address: guoyumin@bjfu.edu.cn

Species distribution models (SDMs) have become an essential tool in ecology, biogeography, evolution, and more recently, in conservation biology. How to generalize species distributions in large undersampled areas, especially with few samples, is a fundamental issue of SDMs. In order to explore this issue, we used the best available presence records for the Hooded Crane (*Grus monacha*, n=33), White-naped Crane (*Grus vipio*, n=40), and Black-necked Crane (*Grus nigricollis*, n=75) in China as three case studies, employing four powerful and commonly used machine learning algorithms to map the breeding distributions of the three species: TreeNet (Stochastic Gradient Boosting, Boosted Regression Tree Model), Random Forest, CART (Classification and Regression Tree) and Maxent (Maximum Entropy Models) Besides, we developed an ensemble forecast by averaging predicted probability of above four models results. Commonly-used model performance metrics (Area under ROC (AUC) and true skill statistic (TSS)) were employed to evaluate model accuracy. Latest satellite tracking data and compiled literature data were used as two independent testing datasets to confront model predictions. We found Random Forest demonstrated the best performance for the most assessment method, provided a better model fit to the testing data, and achieved better species range maps for each crane species in undersampled areas. Random Forest has been generally available for more than 20 years, and by now, has been known to perform extremely well in ecological predictions. However, while increasingly on the rise its potential is still widely underused in conservation, (spatial) ecological applications and for inference. Our results show that it informs ecological and biogeographical theories as well as being suitable for conservation applications, specifically when the study area is undersampled. This method helps to save model-selection time and effort, and it allows robust and rapid assessments

and decisions for efficient conservation.

1   **Why to choose Random Forest to predict rare species distribution with few samples in large**

2   **undersampled areas? Three Asian crane species models provide supporting evidence**

3                   Chunrong Mi[1,2], Falk Huettmann[3], Yumin Guo[1], Xuesong Han[1] and Lijia Wen[1]

4   [1]College of Nature Conservation, Beijing Forestry University, P.O. Box 159, Beijing 100083,

5   China

6   [2]Key Laboratory of Water Cycle and Related Land Surface Processes, Institute of Geographic

7   Sciences and Natural Resources Research, University of Chinese Academy of Sciences, Beijing

8   100101, China

9   [3]EWHALE Lab, Department of Biology and Wildlife, Institute of Arctic Biology, University of

10   Alaska Fairbanks (UAF), 419 Irving I, P.O. Box 757000, AK 99775, USA

11

12   Corresponding author:

13   Yumin Guo

14   College of Nature Conservation, Beijing Forestry University, P.O. Box 159, Beijing 100083,

15   China

16   guoyumin@bjfu.edu.cn

17

18

19

20

21

22

23

24

## ABSTRACT

26  Species distribution models (SDMs) have become an essential tool in ecology, biogeography,

27  evolution, and more recently, in conservation biology. How to generalize species distributions in

28  large undersampled areas, especially with few samples, is a fundamental issue of SDMs. In order

29  to explore this issue, we used the best available presence records for the Hooded Crane (*Grus*

30  *monacha*, n=33), White-naped Crane (*Grus vipio*, n=40), and Black-necked Crane (*Grus*

31  *nigricollis*, n=75) in China as three case studies, employing four powerful and commonly used

32  machine learning algorithms to map the breeding distributions of the three species: TreeNet

33  (Stochastic Gradient Boosting, Boosted Regression Tree Model), Random Forest, CART

34  (Classification and Regression Tree) and Maxent (Maximum Entropy Models) Besides, we

35  developed an ensemble forecast by averaging predicted probability of above four models results.

36  Commonly-used model performance metrics (Area under ROC (AUC) and true skill statistic

37  (TSS)) were employed to evaluate model accuracy. Latest satellite tracking data and compiled

38  literature data were used as two independent testing datasets to confront model predictions. We

39  found Random Forest demonstrated the best performance for the most assessment method,

40  provided a better model fit to the testing data, and achieved better species range maps for each

41  crane species in undersampled areas. Random Forest has been generally available for more than

42  20 years, and by now, has been known to perform extremely well in ecological predictions.

43  However, while increasingly on the rise its potential is still widely underused in conservation,

44  (spatial) ecological applications and for inference. Our results show that it informs ecological and

45  biogeographical theories as well as being suitable for conservation applications, specifically when

46  the study area is undersampled. This method helps to save model-selection time and effort, and it

47    allows robust and rapid assessments and decisions for efficient conservation.

48    *Keywords*: Species distribution models (SDMs), Random Forest, Generality (transferability), Rare

49    species, Undersampled areas, Hooded Crane (*Grus monacha*), White-naped Crane (*Grus vipio*),

50    Black-necked Crane (*Grus nigricollis*)

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

## INTRODUCTION

75   Species distribution models (SDMs) are empirical ecological models that relate species

76   observations to environmental predictors (Guisan & Zimmermann, 2000, Drew et al., 2011).

77   SDMs have become an increasingly important and now essential tool in ecology, biogeography,

78   evolution and, more recently, in conservation biology (Guisan et al., 2013), management

79   (Cushman & Huettmann, 2010), impact assessments (Humphries & Huettmann, 2014) and climate

80   change research (Lei et al., 2011). To generalize and infer from a model, or model transferability

81   is defined as geographical or temporal cross-applicability of models (Thomas & Bovee 1993;

82   Kleyer 2002; Randin et al., 2006). It is one important feature in SDMs, a base-requirement in

83   several ecological and conservation biological applications (Heikkinen et al., 2012). In this study,

84   we used generality (transferability) as the concept of generalizing distribution from sampled areas

85   to unsampled areas (extrapolation beyond the data) in one study area.

86      Detailed distribution data for rare species in large areas are rarely available in SDMs (Pearson

87   et al., 2007; Booms et al., 2010). However, they are the most needed for their conservation to be

88   effective. Collecting and assembling distribution data for species, especially for rare or endangered

89   species in remote wilderness areas is often a very difficult task, requiring a large amount of human,

90   time and funding source (Gwena et al., 2010; Ohse et al., 2009).

91      Recent studies have suggested that machine-learning (ML) methodology, may perform better

92   than the traditional regression-based algorithms (Elith et al., 2006). TreeNet (boosting; Friedman

93    2002), Random Forest (bagging; Breiman, 2001), CART (Breiman et al., 1984) and Maxent

94    (Phillips et al., 2004) are considered to be among the most powerful machine learning algorithms

95    and for common usages (Elith et al., 2006; Wisz et al., 2008; Williams et al., 2009; Lei et al., 2011)

96    and for obtaining powerful ensemble models (Araújo and New 2007; Hardy et al., 2011). Although

97    Heikkinen et al. (2012) compared the four SDMs techniques' transferability in their study, they

98    did not test with rare species and few samples in undersampled areas. It is important to understand

99    that the software platform of the former three algorithms (Boosted Regression Trees, Random

100   Forest and CARTs) applied by Heikkinen et al. (2012) from the R software ("BIOMOD"

101   framewok) comes without a GUI and lacks sophisticated optimization and fine-tuning, but as they

102   are commonly used though by numerous SDM modelers. Instead, we here run these models in the

103   Salford Predictive Modeler (SPM) by Salford Systems Ltd. These algorithms in SPM are further

104   optimized and improved by one of the algorithm's original co-authors (especially for TreeNet and

105   Random Forest). It runs with a convenient GUI, and produces a number of descriptive results and

106   graphics which are not available in the R version. While this is a commercial software, it is usually

107   available on a 30 days trial version (which suffices for most model runs we know. As well, some

108   of the features of the randomForest R package, most notably the ability to produce partial

109   dependence plots (Herrick 2013), are not directly implemented yet in SPM7 (but they can

110   essentially be obtained by running TreeNet in a Random Forest model).

111      Model generality (transferability) testing could offer particularly powerful for model

112   evaluation (Randin et al., 2006). Independent observations from training data set has been

113   recommended as a more proper evaluations of models (Fielding & Bell 1997; Guisan and

114   Zimmermann 2000). So the use of an independent geographically (Fielding & Haworth, 1995) or

115   temporally (Boyce et al., 2002; Araujo et al., 2005b) testing data set is encouraged to assess the

116   generality of different SDMs techniques. Data from museum specimen, published literature

117   (Graham et al., 2004) as well as tracking are good source to assess model generality

118   (transferability) performance. In addition, how the distribution map links with reality data,

119   especially in undersampled areas where modelers want to make predictions should definitely be

120   as a metric to assess model performance and generalization. Arguably, if model predictions

121   perform very well there, great progress is provided. Whereas, predictions on existing knowledge

122   and data offers less progress. The model prediction and conservation frontier obviously sits in the

123   unknown.

124        In this study, we modeled the best-available data for three species in East Asia as test cases:

125   Hooded Cranes (*Grus monacha*, n=33), White-naped Cranes (*Grus vipio*, n=40) and Black-necked

126   Cranes (*Grus nigricollis*, n=75). Four machine-learning models (TreeNet, Random Forest, CART

127   and Maxent) were applied to map breeding distributions for these three crane species which

128   otherwise lack empirically derived distribution information. In addition, two kinds of independent

129   testing data sets (latest satellite tracking data, and compiled literature data (Threatened Birds of

130   Asia: Collar *et al.,* 2001) were obtained to test the transferability of the four model algorithms.

131   The purpose of this investigation is to explore whether there is a SDM technique among the four

132   algorithms that could generate reliable and accurate distributions with high generality for rare

133   species using few samples but in large undersampled areas? Results from this research could be

134   useful for the detection of rare species and enhance fieldwork sampling in large undersampled

135   areas which would save money and effort, as well as the conservation management of those

136   species.

## MATERIALS AND METHODS

### Species data

In our 13 combined years of field work, we have collected 33 Hooded Crane nests (2002-2014),

40 White-naped Crane nests (2009-2014)，and 75 Black-necked Crane nests (2014) (see Fig. 1),

during breeding seasons. We used these field samples (nests) to represent species presence points
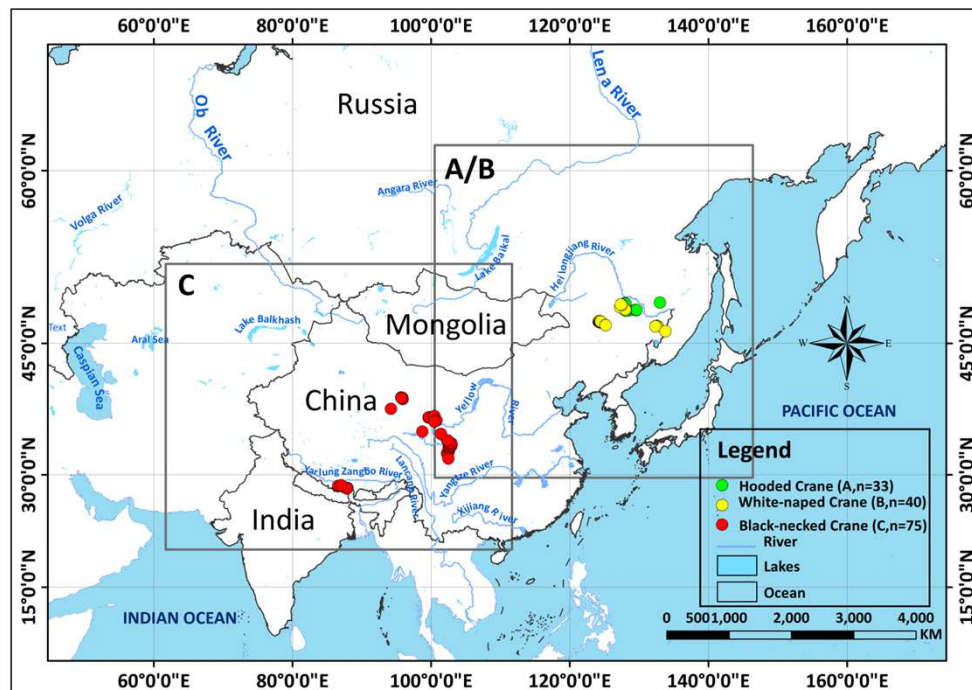
referenced in time and space.

Put Fig. 1 here



Figure 1 Study areas for three species cranes.

### Environmental variables

We used 21 environmental layers at a 30-s resolution in GIS format and that were known to

correlate with bird distribution and as proxies of habitats predictors. They included bio-climatic

factors (bio_1-7, bio_12-15), topographical factors (altitude, slope, and aspect), water factors

150     (distance to river, distance to lake, and distance to coastline), inference factors (distance to road,

151     distance to rail road, and distance to settlements), and land cover factors (for detailed information,

152     see Table 1). Most of these factors were obtained from open access sources. Bio-climate factors

153     were obtained from the WorldClim database, while aspect and slope layer were derived from the

154     altitude layer in ArcGIS, which was also initially obtained from the WorldClim database. Road,

155     railroad, river, lake and coastline and settlement maps were obtained from the Natural Earth

156     database. The land cover map was obtained from the ESA database. We also made models with

157     all 19 bio-climate variables and 10 other environmental variables, and then reduced predictors by

158     AIC, BIC, varclust, PCA and FA analysis. When we compared the distribution maps overlaying

159     with independent data set generated by Random Forest model, we found the model based on 21

160     predictors have the best performance for Hooded Cranes, and the best level for White-naped Crane

161     and Black-necked Cranes (see Supplement S1). Therefore, we decided to constructed models with

162     21 predictors for the all three cranes and four machine-learning techniques. All spatial layers of

163     these environmental variables were resampled to a resolution of 30-s to correspond to that of the

164     bioclimatic variables and for a meaningful high-resolution management scale.

165     <span style="color:red">Put Table 1 here</span>

166     Table 1 Environmental GIS layers used to predict breeding distributions of three cranes.

| Environmental Layers | Description | Source | Website |
|---|---|---|---|
| Bio_1 | Annual mean Temperature (°C) | WorldClim | http://www.worldclim.org/ |
| Bio_2 | Monthly mean (max temp - min temp) (°C) | WorldClim | http://www.worldclim.org/ |

| | | | |
|---|---|---|---|
| Bio_3 | Isothermality (BIO2/BIO7) (*100℃) | WorldClim | http://www.worldclim.org/ |
| Bio_4 | Temperature seasonality (standard deviation *100℃) | WorldClim | http://www.worldclim.org/ |
| Bio_5 | Max temperature of warmest month (℃) | WorldClim | http://www.worldclim.org/ |
| Bio_6 | Min temperature of Coldest month (℃) | WorldClim | http://www.worldclim.org/ |
| Bio_7 | Annual temperature range (BIO5-BIO6) (℃) | WorldClim | http://www.worldclim.org/ |
| Bio_12 | Annual precipitation (mm) | WorldClim | http://www.worldclim.org/ |
| Bio_13 | Precipitation of wettest month (mm) | WorldClim | http://www.worldclim.org/ |
| Bio_14 | Precipitation of driest month (mm) | WorldClim | http://www.worldclim.org/ |
| Bio_15 | Precipitation seasonality (mm) | WorldClim | http://www.worldclim.org/ |
| Altitude | Altitude (m) | WorldClim | http://www.worldclim.org/ |
| Aspect | Aspect (°) | Derived from Altitude | http://www.worldclim.org/ |
| Slope | Slope | Derived from Altitude | http://www.worldclim.org/ |
| Landcover | Land cover | ESA | http://www.esa-landcover-cci.org/ |
| Disroad | Distance to roads (m) | Road layer from Natural Earth | http://www.naturalearthdata.com/ |
| Disrard | Distance to railways (m) | Railroad | http://www.naturalearthdata.com/ |

| | | layer from Natural Earth | |
|---|---|---|---|
| Disriver | Distance to rivers (m) | River layer from Natural Earth | http://www.naturalearthdata.com/ |
| Dislake | Distance to lakes (m) | Lake layer from Natural Earth | http://www.naturalearthdata.com/ |
| Discoastline | Distance to coastline (m) | Coastline layer from Natural Earth | http://www.naturalearthdata.com/ |
| Dissettle | Distance to settlements (m) | Settle layer from Natural Earth | http://www.naturalearthdata.com/ |

## Model development

168    We created TreeNet, Random Forest, CART, Maxent models and ensemble model (averaged

169    value of the former four model results) for Hooded Cranes, White-naped Cranes and Black-naped

170    Cranes. These four model algorithms are considered to be among the most accurate machine

171    learning methods (more information about these four models can be seen in the references by

172    Breiman et al., 1984, Breiman 2001, Friedman 2002, Phillips et al., 2004, Hegel et al., 2010). The

173    first three machine learning models are binary (presence-pseudo absence) models and were

174    handled in Salford Predictive Modeler 7.0 (SPM). For more details on TreeNet, Random Forest

175    and CART in SPM, we refer readers to the user guide document online (https://www.salford-

176    systems.com/products/spm/userguide). Several implementations of these algorithms exist.

177    Approximately 10,000 'pseudo-absence' locations were selected by random sampling across the

178    study area for each species using the freely available Geospatial Modeling Environment (GME;

179 Hawth's Tools; Beyer 2013; see Booms et al., 2010 and Ohse et al., 2009 for examples). We

180 extracted the habitat information from the environmental layers for presence and pseudo-absence

181 points for each crane, and then constructed models in SPM with these data. In addition, we used

182 balanced class weights, and 1000 trees were built for all models to find an optimum within, others

183 used default settings.

184     For the predictions, we created a 'lattice' (equally spaced points across the study area;

185 approximately 5×5 km spacing for the study area). For the lattice, we extracted information from

186 the same environmental layers (Table 1) as described above for each point and then used the model

187 to predict ('score') bird presence for each of the regular lattice points. For visualization, we

188 imported the dataset of spatially referenced predictions ('score file') into GIS as a raster file and

189 interpolated for visual purposes between the regular points using inverse distance weighting (IDW)

190 to obtain a smoothed predictive map of all pixels for the breeding distributions of the three cranes

191 (as performed in Booms et al., 2010 and Ohse et al., 2009). The fourth algorithm we employed,

192 Maxent, is commonly referred to as a presence-only model; we used Maxent 3.3.3k (it can be

193 downloaded for free from http://www.cs.princeton.edu/~schapire/maxent/) to construct our

194 models. To run Maxent, we followed the 3.3.3e tutorial for ArcGIS 10 (Young et al., 2011) and

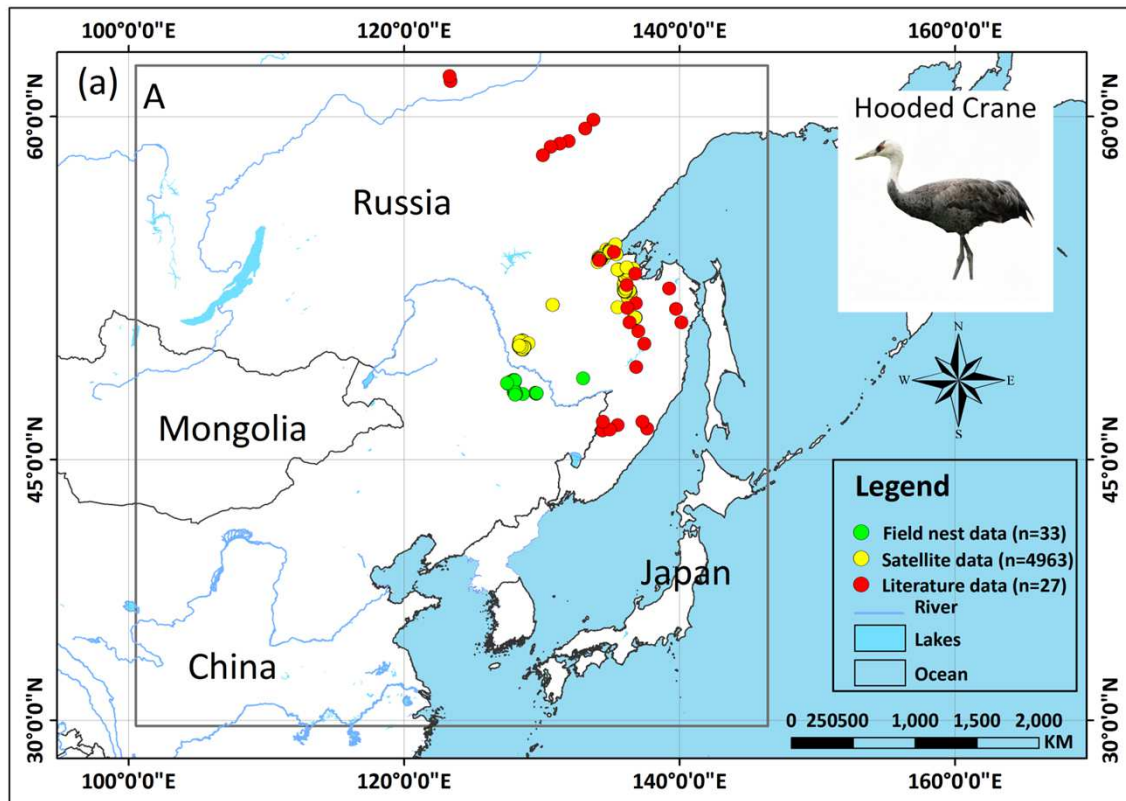195 used default settings.

196 ## Testing data and model assessment

197     We applied two types of testing data in this study: one consisted of satellite tracking data, and

198 the other was represented by data from the literature. Satellite tracking data were obtained from 4

199 individual Hooded Cranes and 8 White-naped Cranes that were tracked in the breeding regions at

200 stopover sites (for more details regarding the information for tracked cranes, please see

201 Supplement S2). The satellite tracking devices could provide 24 data points per day (Databases
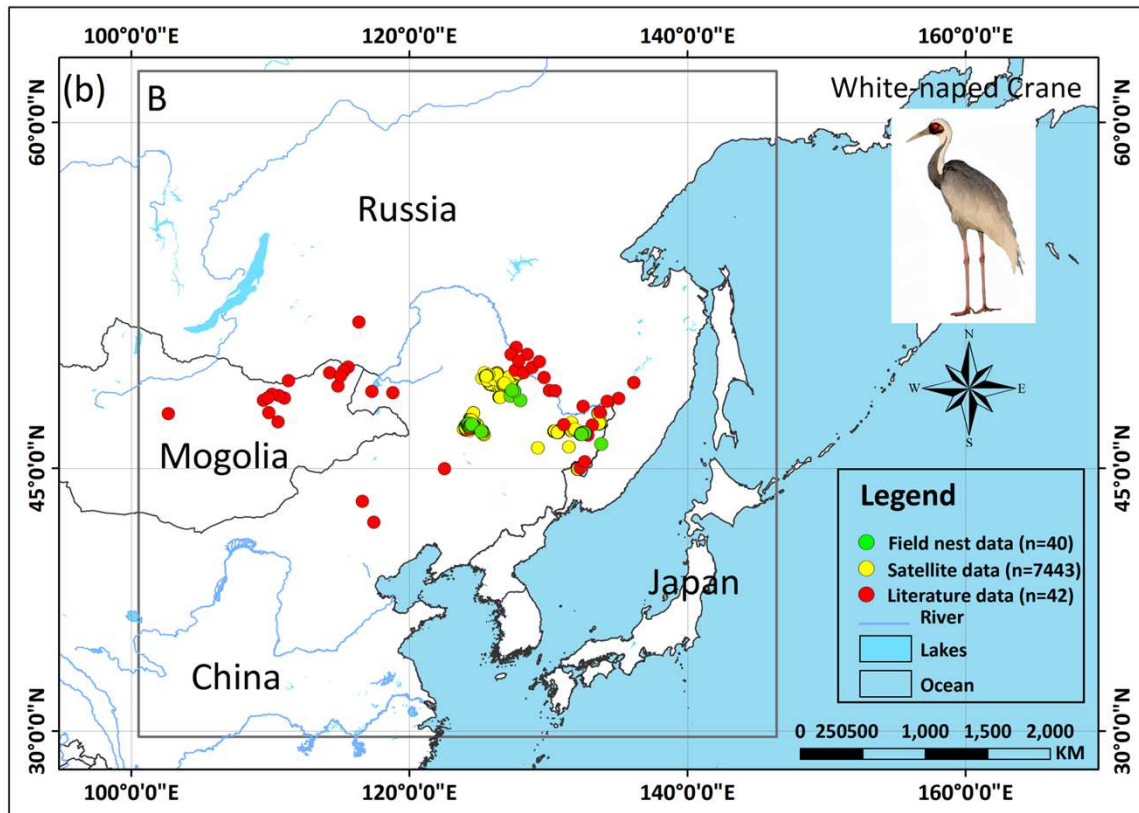
202 could be available upon request). Here, we chose points that had a speed of less than 5 km/h during

203 the period from 1st May to 31th June for Hooded Cranes and 15th April to 15th June for White-naped

204 Cranes as the locations of the breeding grounds for these two cranes. The total numbers of tracking

205 data points were 4,963 and 7,712 (Hooded Cranes and White-naped Crane, respectively. We didn't

206 track Black-necked Cranes, so there was no tracking testing data for this species). The literature

207 data for this study were obtained by geo-referencing the location points of detections from 1980-

208 2000 (ArcGIS 10.1) from Threatened Birds of Asia: the BirdLife International Red Data Book

209 (Collar et al., 2001). From this hardcopy data source, we were able to obtain and digitize 27

210 breeding records for Hooded Cranes, 43 breeding records for White-naped Cranes, and 53 breeding

211 records for Black-necked Cranes (see Fig. 2a, 2b, 2c). We digitized the only crane data for these

212 three species in East-Asia into a database.

213    In addition, we generated 3,000 random points for Hooded Cranes and White-naped Cranes,

214 and 5,000 random points for Black-necked Cranes as testing absence points in their respective

215 study areas. And then, the literature locations (additional presence points for testing) and random

216 points location (testing absence points) that contrasted with the associated predictive value of RIO

217 extracted from the relative prediction map, which were used to calculate receiver operating

218 characteristic (ROC) curves and the true skill statistic (TSS) (Hijmans and Graham, 2006). The

219 area under the ROC curve (AUC) is commonly used to evaluate models in species distributional

220 modeling (Manel *et al.,* 2001, McPherson *et al.,* 2004). TSS was also used to evaluate model

221 performance; we used TSS because it has been increasingly applied as a simple but robust and

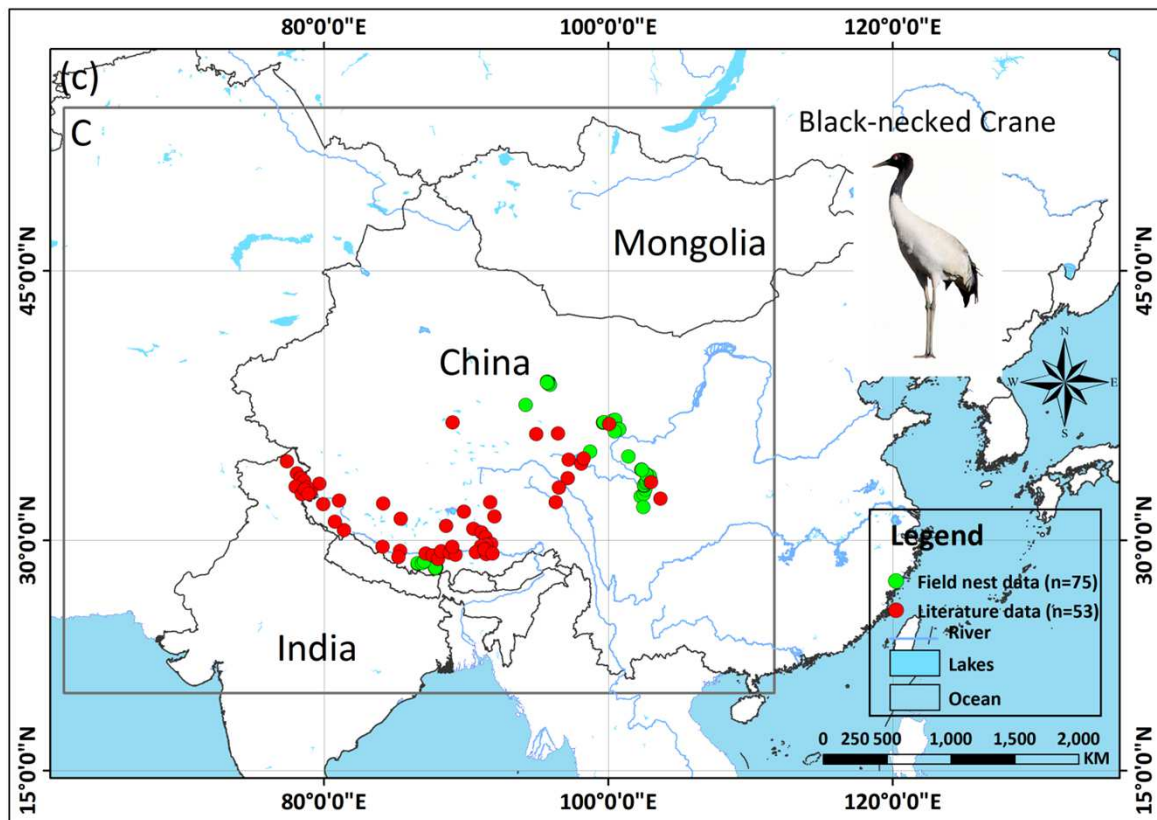222 intuitive measure of the performance of species distribution models (Allouche et al., 2006).

223                                          Put Fig.2 here

224

225

226

227

Figure 2 Detailed study areas showing the presence of and testing data used for the three cranes.

2a) Hooded Cranes, 2b) White-naped Cranes, 2c) Black-necked Cranes.

To assess models transferability, we extracted the predictive value of the relative index of occurrence (RIO) for testing data sets from the prediction maps using GME. We then constructed resulting violin plots for these extracted RIOs to visualize their one-dimensional distribution. This method allowed us to examine the degree of generalizability based on the local area with samples to predict into undersampled areas that are otherwise unsampled in the model development (=areas without training data). In addition, AUC is also commonly used to assess model transferability in our study referring Randin et al. (2006).

# RESULTS

## Model performance

The results for AUC and TSS, two metrics commonly used to evaluate model accuracy, are listed in Table 2. For the four SDMs technique, our results showed that the AUC values for Random Forest were always highest (>0.625), ranking this model in first place, followed by Maxent (>0.558), and then either CART or TreeNet (>=0.500). TSS showed us consistent results, as was the case for AUC, and Random Forest performed the best (>0.250) followed by Maxent (>0.137) for all three crane species, CART took the third place for Black-necked Cranes, and TreeNet performed better than CART for White-naped Cranes. And the results showed there was a trend that the value of these three metrics increased with an increase of nest site samples (33 to 75, Hooded Crane to Black-necked Crane, see Table. 2). Comparing the results of Random Forest with ensemble model, we found their performance were close. Random Forest obtained better model for Hooded Cranes and White-naped Cranes cases, ensemble model performed better for Black-necked Cranes.

<span style="color:red; text-align:center">Put Table 2 here</span>

Table 2 AUC and TSS values for four machine learning models and their ensemble model with three crane species based on literature testing data.
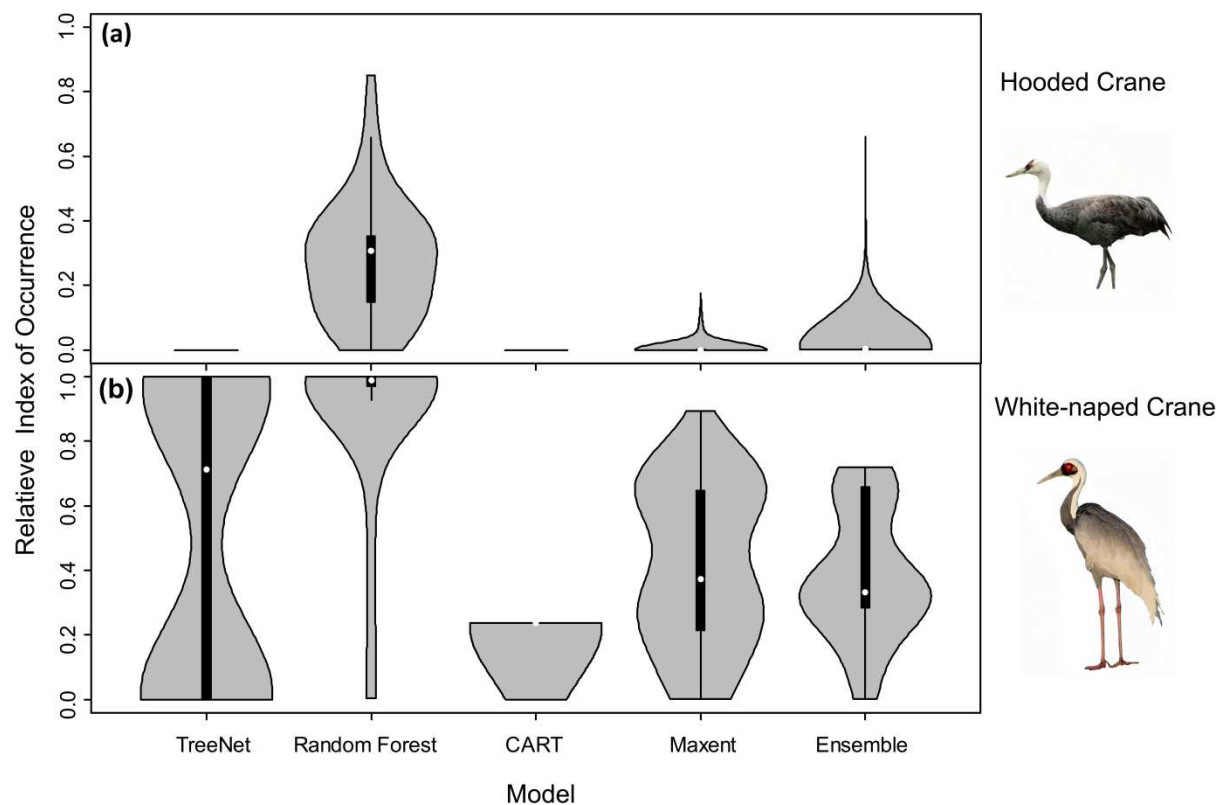
| Accuracy metric (samples) | Species distribution model | | | | |
|---|---|---|---|---|---|
| | TreeNet | Random Forest | CART | Maxent | Ensemble |
| Hooded Crane (*Grus monacha,* n=33 sites) | | | | | |
| AUC | 0.504 | **0.625** | 0.500 | 0.558 | 0.558 |

| | | | | | |
|---|---|---|---|---|---|
| TSS | 0.000 | **0.250** | 0.000 | 0.137 | 0.117 |
| White-naped Crane (*Grus vipio*, n=40 sites) | | | | | |
| AUC | 0.605 | **0.754** | 0.564 | 0.712 | **0.754** |
| TSS | 0.210 | **0.509** | 0.128 | 0.424 | 0.508 |
| Black-necked Crane (*Grus nigricollis*, n=75 sites) | | | | | |
| AUC | 0.528 | 0.830 | 0.672 | 0.805 | **0.843** |
| TSS | 0.055 | 0.660 | 0.345 | 0.611 | **0.686** |

## Model generalization

Violin plots for RIOs with overlaid satellite tracking data (Fig. 3) showed that Random Forest for Hooded Cranes and White-naped Cranes performed better than the other three models. In the Hooded Crane models (Fig. 3a), the RIO for most satellite tracking data indicated that TreeNet, and CART predicted with a value around 0; Ensemble model demonstrated a slightly higher value than the other three models but was still much lower than Random Forest. Fig. 3b indicates the same situation than found in Fig. 3a: Random Forest still performed better than the other three models (median values in Random Forests were close to 1.00). TreeNet had a median RIO value of approximately 0.71, followed by Maxent (median was 0.37) and then ensemble and CART. While some tracking points had a low RIO value in TreeNet, the majority of RIO values for CART remained in the 0.20 range.
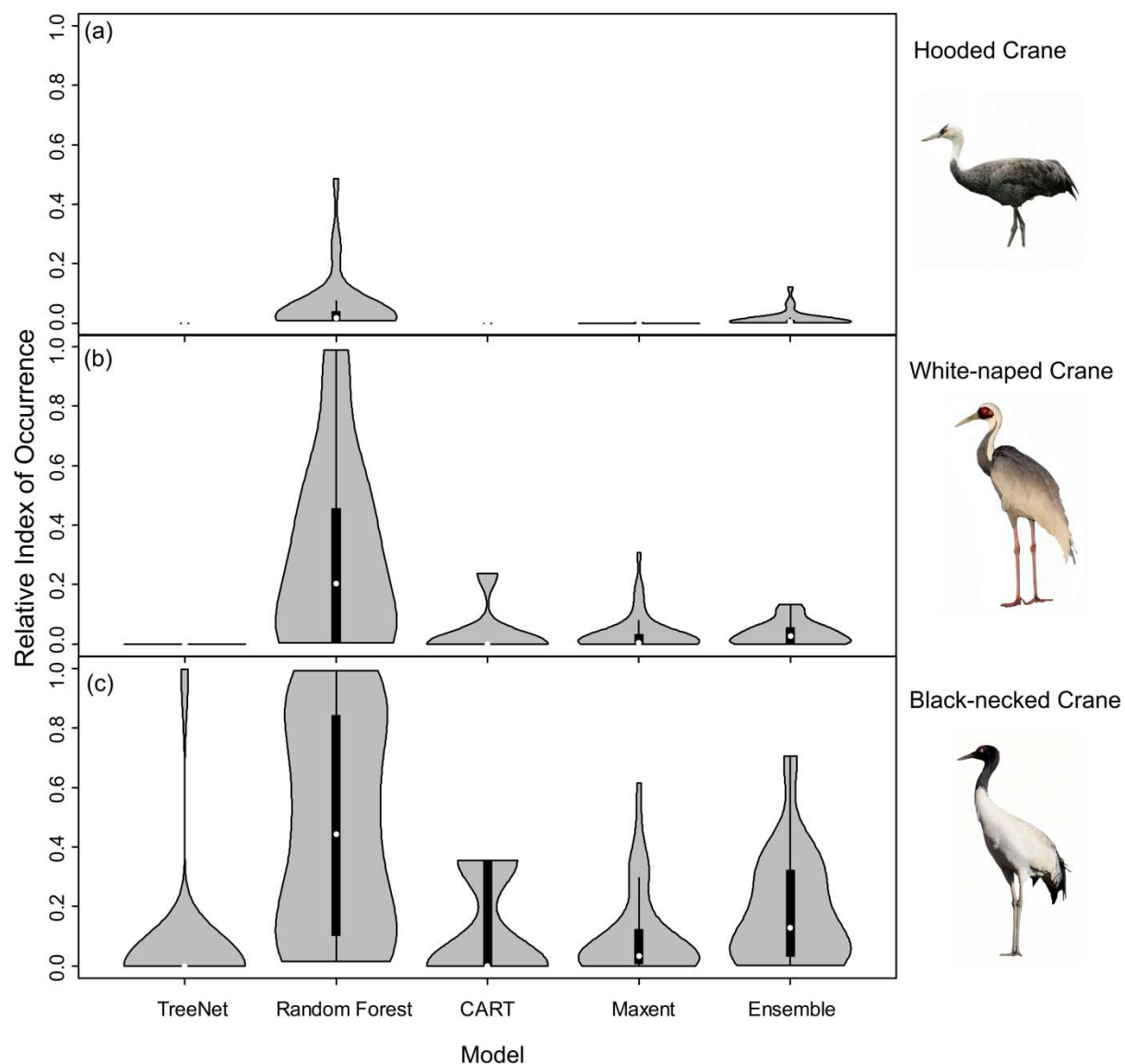
Put Fig. 3 here

266

Figure 3 Violin plots of the Relative Index of Occurrence (RIO) for four SDMs and ensemble

model for Hooded Cranes and White-naped Cranes based on satellite tracking data. 3a) violin plots

of Hooded Cranes, 3b) violin plots of White-naped Cranes.

Violin plots of the RIOs values for the three cranes extracted for the literature data from the

prediction maps (Fig. 4) demonstrated consistent trends (Fig. 3), indicating that Random Forest

performed best across all models of the three species. In Fig. 4a, the RIO values for Random Forest

ranged from 0 to 0.48, and most RIO values were below 0.1; the RIO values for the other three

SDMs method were 0, the ensemble model performed a little bit better. As showed in Fig. 4b, most

RIO values for Random Forest were below 0.7, and the median value was approximately 0.20,

followed by Maxent and then CART. The violin plots for Black-necked Cranes (Fig. 4c) indicated

that TreeNet performed the worst, although there were some pixels that had high RIO values,

278   followed by ensemble and then Maxent. The best performer was still Random Forest, and its RIOs

279   were distributed evenly to a certain extent with a median value of 0.44. The results of AUC, as

280   mentioned in "Model performance" part (Table 2), showed consistent results with violin plots,

281   Random Forest always get the highest value and has the best generalization.

282                                      Put figure 4 here



283

284   Figure 4 Violin plots of Relative Index of Occurrence (RIO) values for four SDMs and ensemble

285   model for three cranes based on calibration data from Threatened Birds of Asia. 4a) Violin plots
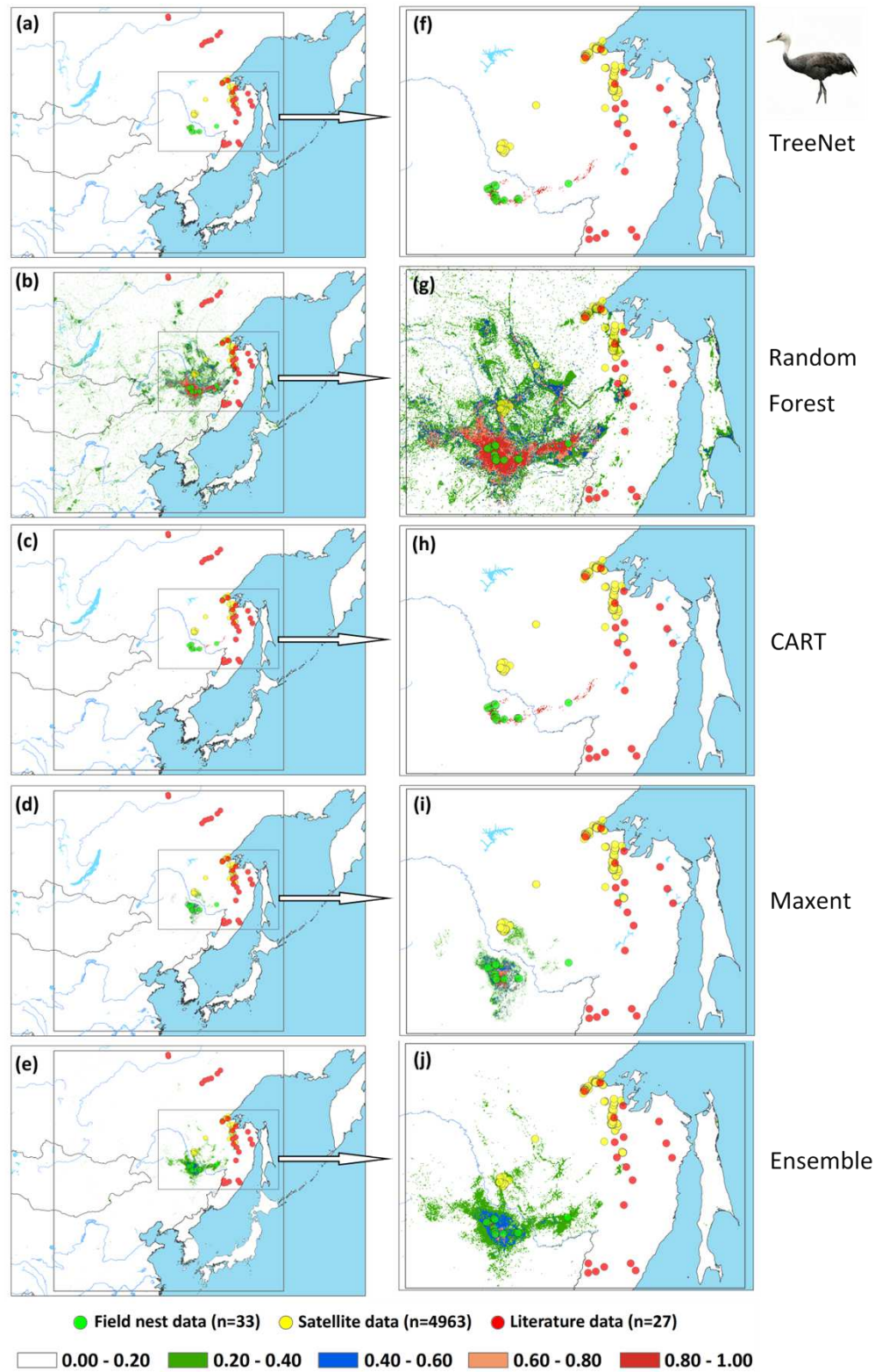
286 for Hooded Cranes, 4b) violin plots for White-naped Cranes, 4c) violin plots for Black-necked

287 Cranes.

## Spatial assessment using a testing data overlay prediction map

289 An assessment of niche prediction beyond the local area where samples were located represents

290 a real test of the generalizability of the model predictions in undersampled areas. This approach

291 was used to evaluate whether testing data (satellite tracking data and literature data) locations

292 matched predictions of the potential distribution area, as a spatial assessment of model

293 performance. It's a spatial and visual method to show the transferability of SDMs from sampled

294 to unsampled areas. From the results (Fig.s 5, 6 and 7. Digital version for each subgraph could be

295 available request), we found that Random Forest demonstrated the strongest performance to handle

296 generality (transferability), and a high fraction of testing data locations were predicted in the

297 distribution areas of the three cranes (Fig.s 5b, 5g, 6b, 6g, 7b, 7g). The order of the generality of

298 the remaining four models was: ensemble model followed by Maxent, CART and then TreeNet.

299 Note, however, that the capacities of these models to predict well in undersampled areas were

300 weaker than Random Forest, it holds particularly for areas that were further away from the sample

301 areas (Fig.s 5, 6 and 7). In addition, we found that the generality increased with sample size (33 to

302 75, Hooded Crane to Black-necked Crane, see Fig.s 5, 6 and 7). This means a higher sample size

303 make models more robust and better to generalize from.

304 <span style="color:red">Put Fig. 5 Here</span>

305

Field nest data (n=33)   Satellite data (n=4963)   Literature data (n=27)

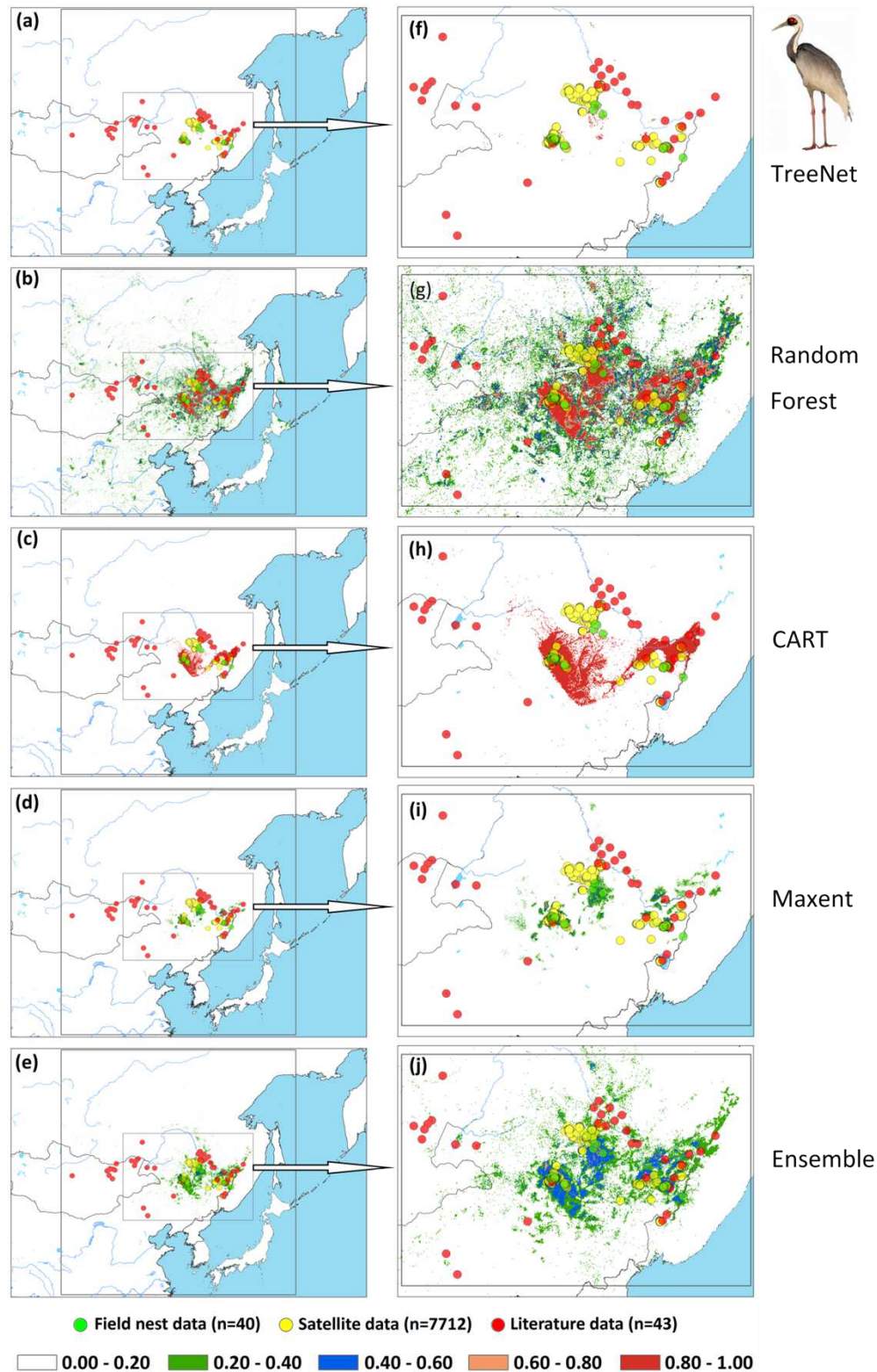0.00 - 0.20   0.20 - 0.40   0.40 - 0.60   0.60 - 0.80   0.80 - 1.00

306

307   Figure 5 Prediction maps for Hooded Cranes and zoomed-in maps showing the four models (TreeNet,

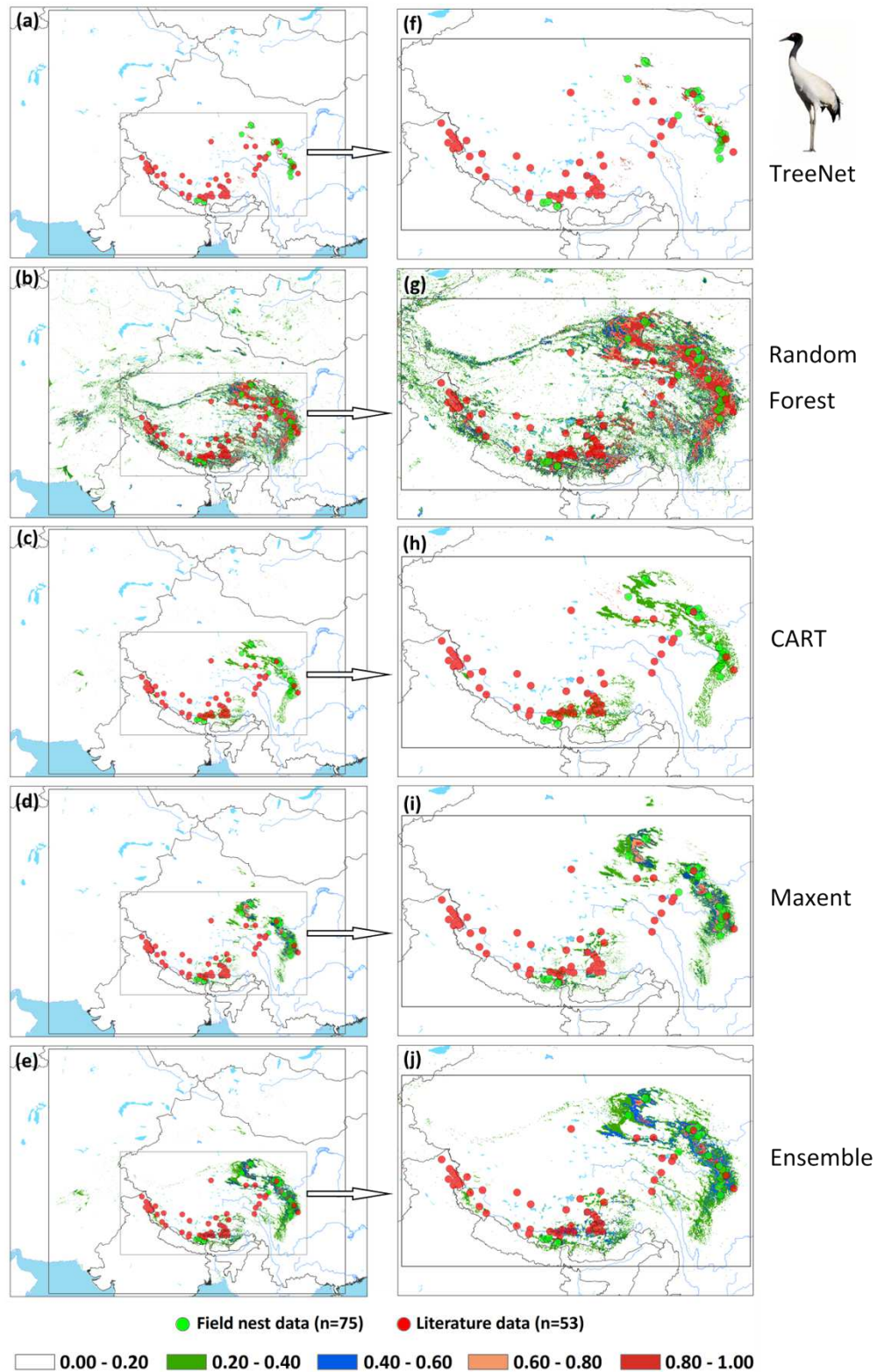308    Random Forest, CART and Maxent) and ensemble model in detail. 5a-5e) prediction map for Hooded

309    Cranes, 5f-5j) zoomed-in map for Hooded Cranes.

Field nest data (n=40)　Satellite data (n=7712)　Literature data (n=43)

0.00 - 0.20　0.20 - 0.40　0.40 - 0.60　0.60 - 0.80　0.80 - 1.00

310

311　Figure 6 Prediction maps for White-naped Cranes and zoomed-in maps showing the four models (TreeNet,

312      Random Forest, CART and Maxent) and ensemble model in detail. 6a-6e) prediction map for White-naped

313      Cranes, 6f-6j) zoomed-in map for White-naped Cranes. Put Fig. 6 Here

Field nest data (n=75)   Literature data (n=53)

0.00 - 0.20   0.20 - 0.40   0.40 - 0.60   0.60 - 0.80   0.80 - 1.00

314

315   Figure 7 Prediction maps for Black-necked Cranes and zoomed-in maps showing the four models (TreeNet,

316  Random Forest, CART and Maxent) and ensemble model in detail. 7a-7e) prediction map for Black-necked

317  Cranes, 7f-7j) zoomed-in map for Black-necked Cranes.

## DISCUSSION

## Model generality (transferability)

320  Estimating species distributions in undersampled areas is a fundamental problem in ecology,

321  biogeography, biodiversity conservation and natural resource management (Drew et al., 2011).

322  That is specifically true for rare and difficult to be detected species and which are usually high on

323  the conservation priority. The use of SDMs has become the method for deriving such estimates

324  (Guisan & Thuiller, 2005; Drew et al., 2011; Guisan et al., 2013) and could contribute to detect

325  new populations of rare species. However, the application of a few samples to project a distribution

326  area widely beyond the sample range is a greater challenge and has rarely been attempted in the

327  literature. And only recently have conservationists realized its substantial value for pro-active

328  decision making in conservation management (see work by Ohse et al., 2010; Drew et al., 2011;

329  Kandel et al., 2015 etc.). Our results based on AUC, violin plots for RIOs and spatial assessment

330  of testing data (satellite tracking data and literature data) all suggest there are difference in the

331  generalization performance of different modeling techniques (TreeNet, Random Forest , CART

332  and Maxent).

333  Moreover, among the acknowledged four rather powerful and commonly used machne-learning

334  techniques, Random Forest (bagging) in SPM usually had the best performance in each case. Our

335  results are in agreement with those of Prasad et al. (2006), Cutler et al. (2007) and Syphard and

336  Franklin (2009) indicating a superiority of Random Forest in such applications. However, initially

337  it appears to run counter to the conclusions off recent paper (Heikkinen et al., 2012) with the poor

338  transferability of Random Forest. But we propose this is due to the fact that many Random Forest

339  implementations exist (see the 100 classifier paper Fernández-Delgado et al., 2014).

340  Here we applied Random Forest in SPM which has been optimized under one of the algorithm's

341  original co-authors, while Heikkinen et al. (2012) just run a basic Random Forest with BIOMOD

342  framework in the R sofeware. The differences are known to be rather big (see Herrick 2013).

343  Furthermore, Maxent, a widely used SDM method enjoyed by many modelers (Phillips et al.,

344  2006; Peterson et al. 2007; Phillips and Dudík 2008; Li et al., 2015, etc.), didn't perform so good

345  in regards to transferability in this study. This contrasts to those of Elith et al. (2006) and Heikkinen

346  et al. (2012), where Manxent and GBM perform well. We infer this may be caused by sample size

347  used as training data. When the sample size increased (33 to 75), the AUC and TSS value of all

348  models rose (Table 2). This indicates that higher sample sizes make models more robust and

349  performing better. Sample sizes of 33 presence points still favor by Random Forest.

350  In Random Forest, random samples from rows and variables are used to build hundreds of trees.

351  Each individual tree is constructed from a bootstrap sample and split at each node by the best

352  predictor from a very small, randomly chosen subset of the predictor variable pool (Herrick, 2013).

353  These trees comprising the forest are each grown to maximal depth, and predictions are made by

354  averaged trees through 'voting' (Breiman et al., 2006). This algorithm avoiding overfitting by

355  controlling the number of predictors randomly used at each split, using means of out-of-bag (OOB)

356  samples to calculate an unbiased error rate. And also, Random Forest in SPM utilizes additional

357  specific fine-tuning for best performance.

## RIOs of random points

359  In order to explore whether Random Forest created higher RIOs for prediction maps in each grid,

360  which would result higher RIOs of testing data, we generated 3,000 random points for Hooded

361  Cranes and White-naped Cranes, 5000 random points for Black-necked Cranes in their related

362   projected study areas. We made violin plots for RIOs of random points (Fig. 8), and found that

363   more RIO values of random points for Maxent, Random Forest and ensemble models were close

364   to the lower value, and then followed by TreeNet. The distribution shapes of Random Forest,

365   Maxent and ensemble model are more similar to the real distribution of species in the real world.

366   The RIOs of White-naped Crane extracted from the CART model distributed in the range of the

367   low value. That means there were no points located in the high RIO areas of cranes, and which is

368   unrealistic. Consequently, we argued that Random Forest did not create higher RIOs for prediction

369   maps in each grid in our study.

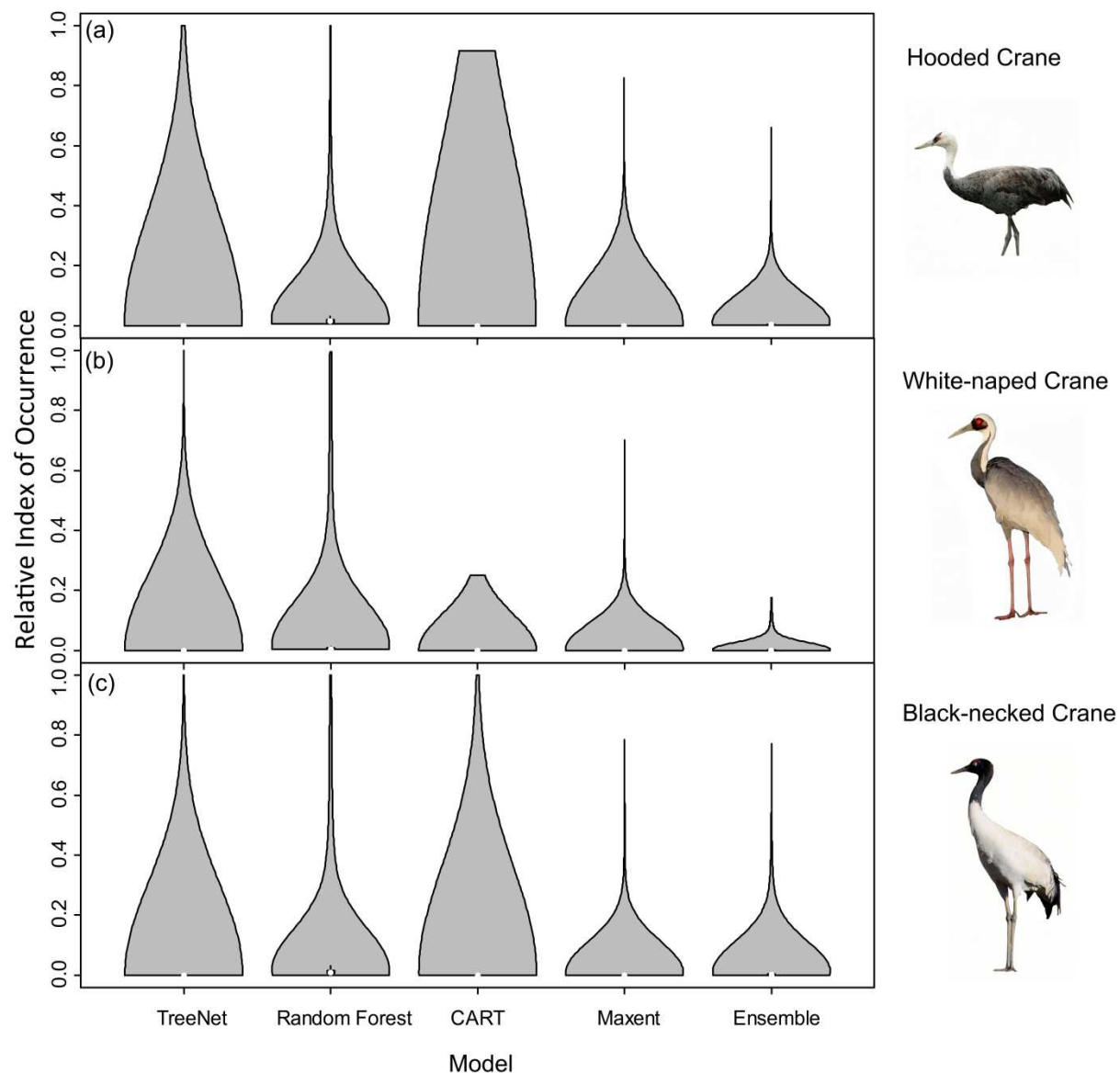370                                    Put Fig. 8 here

371

Figure 8 Violin plots of Relative Index of Occurrence (RIO) values for four SDMs and ensemble

model for three cranes based on calibration data from Threatened Birds of Asia. 4a) Violin plots

for Hooded Cranes, 4b) violin plots for White-naped Cranes, 4c) violin plots for Black-necked

Cranes.

## Models with small sample sizes

Conservation biologists are often interested in rare species and seek to improve their conservation.

378 These species usually have limited number of available occurrence records, which poses

379 challenges for the creation of accurate species distribution models when compared with models

380 developed with greater numbers of occurrences (Stockwell & Peterson, 2002; McPherson et al.,

381 2004; Hernandez et al., 2006). In this study, we used three crane species as case studies, and their

382 occurrence records (nests) totaled 33, 40, and 75, respectively (considering the small numbers of

383 samples and given that a low fraction of the area was sampled in the large projected area). In our

384 models, we found that model fit (AUC and TSS, see Table 2) of Random Forest that had the highest

385 index, while Maxent usually ranked second. In addition, we found that models with few presence

386 samples can also generate accurate species predictive distributions (Fig. 3 to 7) with the Random

387 Forest method. Of course, models constructed with few samples underlie the threat of being biased

388 more because few samples usually had not enough information including all distribution gradients

389 conditions of a species, especially for places far away from the location of training presence points.

390 However, the potential distribution area predicted by SDMs could become as the place where

391 scholars could look for the birds (additional fieldwork sampling). And also, these places could be

392 used as diffusion or reintroduction areas!

### Evaluation methods

394 In this study, we applied two widely-used assessment methods (AUC and TSS) in SDMs (Table

395 2). For evaluation of these three values we used the approach recommended by Fielding & Bell

396 (1997), and Allouche et al. (2006), we found our model usually didn't obtain perfect performance,

397 and some of them were fair. However, for macro-ecology this more than reasonable and ranks

398 rather high. It's a good conservation progress! We identified Random Forest as always the highest

399 performing. These results are consistent with the results of violin plots of the Relative Index of

400 Occurrence (RIO) using tracking as well as literature data (Fig.s 3, 4), and well as matching the

401   spatial assessment results (Fig.s 5-7). And we recommend when modelers assess model

402   performance they should not only depend alone on some metric (such as AUC and TSS), but also

403   should base their assessments on the combined use of visualization and expert knowledge. That

404   means modelers should also assess how the species distribution map actually looks and how it

405   links with real data (see Huettmann & Gottschalk 2011). Spatial assessment metrics from

406   alternative data should matter the most. Expert experience and ecological common knowledge of

407   the species of interest could sometimes also be highly effective (Drew & Perera, 2011), albeit

408   nonstandard, evaluation methods (see Kandel et al., 2015 for an example). Additionally, one

409   alternative method for rapid assessment we find is to use a reliable SDM, and thus Random Forest

410   may be a good choice in the future given our consistent results (Fig.s 3 to 7, Tables 3 to 5) in this

411   study, which involved three species, a vast landscape to conserve, and only limited data. Our work

412   helps to inform conservation decisions for cranes in Northeast Asia.

413   **Limitations and future work**

414   Our study is not without limitations: 1) so far, only three species of cranes are used as a test case

415   in our study. That's because nest data for rare species in remote areas are usually sparse; 2) all our

416   species study areas are rather vast and confined to East-Asia. For future, we would apply Random

417   Forest in more species and in more geography conditions with different distributed feature for a

418   first rapid assessment and baseline mandatory for better conservation. Then we would apply our

419   prediction results in specifically targeted fieldwork sampling campaigns and assess the model

420   accuracy with field survey results (ground-truthing) and more new satellite tracking data. This is

421   to be fed directly into the conservation management process.

## 422 ACKNOWLEDGEMENTS

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

## REFERENCES

Allouche OA, Tsoar, Kadmon R. 2006. Assessing the accuracy of species distribution models: prevalence, kappa and the true skill statistic (TSS). *Journal of Applied Ecology* 43(6):1223-1232.

Araújo MB, New M. 2007. Ensemble forecasting of species distributions. Trends in Ecology & Evolution 22(1):42-47.

Araújo MB, Whittaker R, Ladle R, Erhard M. 2005. Reducing uncertainty in projections of extinction risk from climate change. Global Ecology & Biogeography 14(6):529-538.

Ashtonw C, Perera AH. 2010. Expert Knowledge as a Basis for Landscape Ecological Predictive Models. Predictive Species & Habitat Modeling in Landscape Ecology:229-248.

Beyer H. 2013. Hawth's Analysis Tools for ArcGIS version 3.27 (software). in.

Booms TL, Huettmann F, Schempf PF. 2010. Gyrfalcon nest distribution in Alaska based on a predictive GIS model. Polar biology 33(3):347-358.

Boyce MS, Vernier PR, Nielsen SE, Schmiegelow FK. 2002. Evaluating resource selection functions. Ecological Modelling 157(2-3):281–300.

Braunisch V, Coppes J, Arlettaz R, Suchant R, Schmid H, Bollmann K. 2013. Selecting from correlated climate variables: a major source of uncertainty for predicting species distributions under climate change. Ecography 36(9):971-983.

Breiman L. 2001. Random forests. Machine learning 45(1):5-32.

Breiman L., Friedman J, Stone CJ, Olshen RA. 1984. Classification and regression trees. CRC press.

Bucklin DN, Basille M, Benscoter AM, Brandt LA, Mazzotti FJ, Romanach SS, Speroterra C, Watling JI. 2015. Comparing species distribution models constructed with different subsets of

466    environmental predictors. Diversity and Distributions 21(1):23-35.

467    Cohen J. 1960. A Coefficient of Agreement for Nominal Scales. Educational and Psychological

468    Measurement 20(1):37-46.

469    Collar NJ, Crosby R, Crosby M. 2001. Threatened birds of Asia: the BirdLife International red

470    data book. Volume 1.BirdLife International Cambridge, UK.

471    Cushman SA, Huettmann F. 2010. Spatial Complexity, Informatics, and Wildlife Conservation.

472    Springer, Springer Tokyo Berlin Heidelberg New York.

473    Cutler DR, Edwards Jr TC, Beard KH, Cutler A, Hess KT, Gibson J, Lawler JJ. 2007. Random

474    forests for classification in ecology. Ecology 88(11):2783-2792.

475    Drew CA, Perera AH. 2011. Expert knowledge as a basis for landscape ecological predictive

476    models. Pages 229-248 in Predictive Species and Habitat Modeling in Landscape Ecology.

477    Springer.

478    Drew CA, Wiersma Y, Huettmann F. 2011. Predictive species and habitat modeling in landscape

479    ecology. Springer.

480    Elith J, Graham CH, Anderson RP, Dudík M, Ferrier S, Guisan A, Hijmans RJ, Huettmann F,

481    Leathwick JR, Lehmann A, Li J, Lohmann LG, Loiselle BA, Manion G, Moritz C, Nakamura

482    M, Nakazawa Y, Overton JMM, Peterson AT, Phillips SJ, Richardson K, Scachetti-Pereira R,

483    Schapire RE, Soberón J, Williams S, Wisz MS, Zimmermann NE. 2006. Novel methods

484    improve prediction of species' distributions from occurrence data. Ecography 29(2):129-151.

485    Eskildsen, A., P. C. Roux, R. K. Heikkinen, T. T. Høye, W. D. Kissling, J. Pöyry, M. S. Wisz, and

486    M. Luoto. 2013. Testing species distribution models across space and time: high latitude

487    butterflies and recent warming. Global ecology and biogeography 22(12):1293-1303.

488    Estes L, Bradley B, Beukes H, Hole D, Lau D, Oppenheimer M, Schulze R, Tadross M, Turner

489     W. 2013. Comparing mechanistic and empirical model projections of crop suitability and

490     productivity: implications for ecological forecasting. Global ecology and biogeography

491     22(8):1007-1018.

492   Fernández-Delgado M, Cernadas E, Barro S, Amorim D. 2014. Do we need hundreds of classifiers

493     to solve real world classification problems? The Journal of Machine Learning Research

494     15(1):3133-3181.

495   Ferrier S, Watson G, Pearce J, Drielsma M. 2002. Extended statistical approaches to modelling

496     spatial pattern in biodiversity in northeast New South Wales. I. Species-level modelling.

497     Biodiversity & Conservation 11(12):2275-2307.

498   Fielding AH, Bell JF. 1997. A review of methods for the assessment of prediction errors in

499     conservation presence/absence models. Environmental conservation 24(1):38-49.

500   Fielding AH, Haworth PF. 1995. Testing the Generality of Bird‐Habitat Models. Conservation

501     biology 9(6):1466-1481.

502   Ferrier S, Watson G, Pearce J, Drielsma M. 2002. Extended statistical approaches to modelling

503     spatial pattern in biodiversity in northeast New South Wales. I. Species-level modelling.

504     Biodiversity & Conservation 11(12):2275-2307.

505   Friedman JH. 2002. Stochastic gradient boosting. Computational Statistics & Data Analysis

506     38(4):367-378.

507   Graham CH, Ferrier S, Huettman F, Moritz C, Peterson AT. 2004. New developments in museum-

508     based informatics and applications in biodiversity analysis. Trends in Ecology & Evolution

509     19(9):497-503.

510   Guillera‐Arroita G, Lahoz‐Monfort JJ, Elith J, Gordon A, Kujala H, Lentini PE, McCarthy MA,

511     Tingley R, Wintle BA. 2015. Is my species distribution model fit for purpose? Matching data

512    and models to applications. Global ecology and biogeography 24(3):276-292.

513    Guisan A, Thuiller W. 2005. Predicting species distribution: offering more than simple habitat

514        models. Ecology letters 8(9):993-1009.

515    Guisan A, Tingley R, Baumgartner JB, Naujokaitis‑Lewis I, Sutcliffe PR, Tulloch AIT, Regan

516        TJ, Brotons L, Mcdonald‑Madden E, Mantyka‑Pringle C. 2013. Predicting species

517        distributions for conservation decisions. Ecology letters 16(12):1424-1435.

518    Guisan A, Zimmermann NE. 2000. Predictive habitat distribution models in ecology. Ecological

519        Modelling 135(2):147-186.

520    Gwena LL, Robin E, Erika F, Guisan A. 2010. Prospective sampling based on model ensembles

521        improves the detection of rare species. Ecography 33(6):1015-1027.

522    Hanley JA, McNeil BJ. 1982. The meaning and use of the area under a receiver operating

523        characteristic (ROC) curve. Radiology 143(1):29-36.

524    Hanley JA, McNeil BJ. 1983. A method of comparing the areas under receiver operating

525        characteristic curves derived from the same cases. Radiology 148(3):839-843.

526    Hardy SM, Lindgren M., Konakanchi H, Huettmann F. 2011. Predicting the distribution and

527        ecological niche of unexploited snow crab (Chionoecetes opilio) populations in Alaskan

528        waters: a first open-access ensemble model. Integrative and comparative biology 51(4):608-

529        622.

530    Hegel TM, SA Cushman, J Evans, Huettmann F. 2010. Current State of the Art for Statistical

531        Modelling of Species Distributions. Spatial Complexity, Informatics, and Wildlife

532        Conservation:273-311.

533    Heikkinen RK, Marmion M, Luoto M. 2012. Does the interpolation accuracy of species

534        distribution models come at the expense of transferability? Ecography 35(3):276-288.

535   Hernandez PA, Graham CH, Master LL, Albert DL. 2006. The effect of sample size and species

536        characteristics on performance of different species distribution modeling methods. Ecography

537        29(5):773-785.

538   Herrick K. 2013. Predictive Modeling of Avian Influenza in Wild Birds. Veterinary Research.

539   Hijmans RJ, Graham CH. 2006. The ability of climate envelope models to predict the effect of

540        climate change on species distributions. Global Change Biology 12(12):2272-2281.

541   Huettmann F, Gottschalk T. 2011. Simplicity, Model Fit, Complexity and Uncertainty in Spatial

542        Prediction Models Applied Over Time: We Are Quite Sure, Aren't We? Pages 189-208 in

543        Predictive Species and Habitat Modeling in Landscape Ecology. Springer.

544   Humphries GRW, Huettmann F. 2014. Putting models to a good use: a rapid assessment of Arctic

545        seabird biodiversity indicates potential conflicts with shipping lanes and human activity.

546        Diversity and Distributions 20(4):478-490.

547   Jiguet F, Barbet-Massin M, Chevallier D. 2011. Predictive distribution models applied to satellite

548        tracks: modelling the western African winter range of European migrant Black Storks Ciconia

549        nigra. Journal of Ornithology 152(1):111-118.

550   Kandel K, Huettmann F, Suwal MK, Regmi GR, Nijman V, Nekaris K, Lama ST, Thapa A,

551        Sharma HP, Subedi TR. 2015. Rapid multi-nation distribution assessment of a charismatic

552        conservation species using open access ensemble model GIS predictions: Red panda (Ailurus

553        fulgens) in the Hindu-Kush Himalaya region. Biological Conservation 181:150-161.

554   Keith DA, Elith J, Simpson CC. 2014. Predicting distribution changes of a mire ecosystem under

555        future climates. Diversity and Distributions 20(4):440-454.

556   Kessler A, Batbayar N, Natsagdorj T, Batsuur D, Smith A. 2013. Satellite telemetry reveals

557        long‑distance migration in the Asian great bustard Otis tarda dybowskii. Journal of Avian

558      Biology 44(4):311-320.

559   Kleyer M. 2002. Validation of plant functional types across two contrasting landscapes. Journal of

560      Vegetation Science 13(2):167-178.

561   Lei, Z., L. Shirong, S. Pengsen, and WangTongli. 2011. Comparative evaluation of multiple

562      models of the effects of climate change on the potential distribution of Pinus massoniana.

563      Chinese Journal of Plant Ecology 35(11):1091-1105.

564   Li, R., M. Xu, M. H. G. Wong, S. Qiu, X. Li, D. Ehrenfeld, and D. Li. 2015. Climate change

565      threatens giant panda protection in the 21st century. Biological Conservation 182:93-101.

566   Maggini R, Lehmann A, Zbinden N, Zimmermann NE, Bolliger J, Schröder B, Foppen R, Schmid

567      H, Beniston M, Jenni L. 2014. Assessing species vulnerability to climate and land use change:

568      the case of the Swiss breeding birds. Diversity and Distributions 20(6):708-719.

569   Manel S, Williams HC, Ormerod SJ. 2001. Evaluating presence–absence models in ecology: the

570      need to account for prevalence. Journal of Applied Ecology 38(5):921-931.

571   McPherson J, Jetz W, Rogers DJ. 2004. The effects of species' range sizes on the accuracy of

572      distribution models: ecological phenomenon or statistical artefact? Journal of Applied Ecology

573      41(5):811-823.

574   Mingchang C, Guangsheng Z, Ensheng W. 2005. Application and comparison of generalized

575      models and classification and regression tree in simulating trees species distribution. ACTA

576      ECOLOGICA SINICA 25(8):2031-2040.

577   Navarro‐Cerrillo R, Hernández‐Bermejo J, Hernández‐Clemente R. 2011. Evaluating models

578      to assess the distribution of Buxus balearica in southern Spain. Applied Vegetation Science

579      14(2):256-267.

580   Ohse B, Huettmann F, Ickert-Bond SM, Juday GP. 2009. Modeling the distribution of white spruce

581   (Picea glauca) for Alaska with high accuracy: an open access role-model for predicting tree

582   species in last remaining wilderness areas. Polar biology 32(12):1717-1729.

583   Pearson RG, Raxworthy CJ, Nakamura M, Peterson AT. 2007. Predicting species distributions

584   from small numbers of occurrence records: a test case using cryptic geckos in Madagascar.

585   Journal of Biogeography 34(1):102-117.

586   Peterson AT., Monica P, Muir E. 2007. Transferability and model evaluation in ecological niche

587   modeling: a comparison of GARP and Maxent. Ecography 30(4):550–560.

588   Phillips SJ, Anderson RP, Schapire RE. 2006. Maximum entropy modeling of species geographic

589   distributions. Ecological Modelling 190(3):231-259.

590   Phillips SJ, Dudík M. 2008. Modeling of species distributions with Maxent: new extensions and a

591   comprehensive evaluation. Ecography 31(2):161-175.

592   Phillips SJ, Dudík M, Schapire RE. A maximum entropy approach to species distribution

593   modeling. ACM, 2004.

594   Prasad AM, Iverson LR, Liaw A. 2006. Newer Classification and Regression Tree Techniques:

595   Bagging and Random Forests for Ecological Prediction. Ecosystems 9(2):181-199.

596   Randin CF, Dirnböck T, Dullinger S, Zimmermann NE, Zappa M, Guisan A. 2006. Are niche-

597   based species distribution models transferable in space? Journal of Biogeography 33(10):1689-

598   1703.

599   Romo H, García-Barros E, Márquez AL, Moreno JC, Real R. 2014. Effects of climate change on

600   the distribution of ecologically interacting species: butterflies and their main food plants in

601   Spain. Ecography 37(11):1063-1072.

602   Stockwell DR, Peterson AT. 2002. Effects of sample size on accuracy of species distribution

603   models. Ecological Modelling 148(1):1-13.

604    Stokes KL, Broderick AC, Canbolat AF, Candan O, Fuller WJ, Glen F, Levy Y, Rees AF, Rilov

605        G, Snape RT, Stott I, Tchernov D, Godley BJ. 2015. Migratory corridors and foraging hotspots:

606        critical habitats identified for Mediterranean green turtles. Diversity and Distributions 21(6):

607        665-674.

608    Swets JA 1988. Measuring the accuracy of diagnostic systems. Science 240(4857):1285-1293.

609    Syphard DA, Franklin J. 2009. Differences in spatial predictions among species distribution

610        modeling methods vary with species traits and environmental predictors. Ecography

611        32(6):907-918.

612    Thomas JA, Bovee KD. 1993. Application and testing of a procedure to evaluate transferability of

613        habitat suitability criteria. Regulated rivers 8:285-285.

614    Thuiller W. 2003. BIOMOD–optimizing predictions of species distributions and projecting

615        potential future shifts under global change. Global Change Biology 9(10):1353-1362.

616    Zhai T, Li X. 2012. Climate change induced potential range shift of the crested ibis based on

617        ensemble models. ACTA ECOLOGICA SINICA 32(8):2361-2370 (in Chinese).

618    Williams JN, Seo C, Thorne J, Nelson JK, Erwin S, O'Brien JM, Schwartz MW. 2009. Using

619        species distribution models to predict new occurrences for rare plants. Diversity and

620        Distributions 15(4):565-576.

621    Wisz MS, Hijmans RJ, Li J, Peterson AT, Graham CH, Guisan A. 2008. Effects of sample size on

622        the performance of species distribution models. Diversity and Distributions 14(5):763-773.

623    Yen P, Huettmann F, Cooke F. 2004. A large-scale model for the at-sea distribution and abundance

624        of Marbled Murrelets (Brachyramphus marmoratus) during the breeding season in coastal

625        British Columbia, Canada. Ecological Modelling 171(4):395-413.

626    Young N, Carter L, Evangelista P. 2011. A MaxEnt Model v3.3.3e Tutorial.

627 Zhang M, Zhou Z, Chen W, Cannon CH, Raes N, Slik JWF. 2014. Major declines of woody plant

628 species ranges under climate change in Yunnan, China. Diversity and Distributions 20(4):405-

629 415.

630

631

632

633

634

635

636