1 **A generative deep learning approach for global species distribution**

2 **prediction**

3

4

5

6 Yujing Yan[1]*, Bin Shao[2]*, Charles C. Davis[1]

7

8

9

10 [1] Department of Organismic and Evolutionary Biology, Harvard University Herbaria,

11 Harvard University, Cambridge MA 02138, USA

12 [2] Department of Molecular and Cellular Biology, Harvard University, Cambridge MA 02138,

13 USA

14

15

16

17 **\*Corresponding authors:**

18 Yujing Yan: yjyan7@gmail.com

19 Bin Shao: shaobinlx@gmail.com

20

21

22

23

24

25

26

27 **Abstract**

28 Anthropogenic pressures on biodiversity necessitate efficient and highly scalable methods

29 to predict global species distributions. Current species distribution models (SDMs) face

30 limitations with large-scale datasets, complex interspecies interactions, and data quality.

31 Here, we introduce EcoVAE, a framework of autoencoder-based generative models trained

32 separately on nearly 124 million georeferenced occurrences from taxa including plants,

33 butterflies and mammals, to predict their global distributions at both genus and species

34 levels. EcoVAE achieves high precision and speed, captures underlying distribution

35 patterns through unsupervised learning, and reveals interspecies interactions via *in silico*

36 perturbation analyses. Additionally, it evaluates global sampling efforts and interpolates

37 distributions without relying on environmental variables, offering new applications for

38 biodiversity exploration and monitoring.

39

40

41

42

43

44

45

46

47

48

49

50

51

2

## Main

Anthropogenic pressures have intensified the need for efficient and scalable methods to predict species distributions for attaining a comprehensive picture of biodiversity. Over the past two decades, species distribution modeling (SDM) has become an essential tool for this purpose [1–3], typically using species occurrence data and environmental variables to predict distributions through statistical and machine learning (ML) algorithms [4–7]. While the mobilization of vast amounts of specimen records and the rapid accumulation of observational data have greatly promoted the development of SDMs [8–10], several challenges still persist.

First, current SDMs struggle to handle large-scale datasets in our big data era [8–12], especially for modeling species assemblages. Traditional methods can only address these tasks through computationally intensive "stacking methods" that have limited scalability. Second, most SDMs overlook complex interspecies interactions, limiting their ecological relevance and utility in modeling community dynamics [13–16]. Third, the accuracy of current techniques depends heavily on input data quality and resolution [17,18]. Although platforms like GBIF and eBird provide valuable datasets for exploration, these are often biased by varying observation efforts across taxonomic groups and regions [19–21]. Finally, reliance on environmental variables introduces additional issues, including collinearity and limited availability in certain regions, further constraining model accuracy and applicability[22].

Generative models are a type of deep learning model that captures complex and nonlinear relationships between input variables. They have been widely adopted in various fields, including natural language processing [23], image generation [24], data capturing [25], and biomedicine [26]. Among these models, autoencoders are designed to compress and reconstruct data in an unsupervised way, making it effective for data denoising, interpolation, and handling randomly missing data [27].

Here, we present for the first time an autoencoder-based framework, named Ecological Variational Autoencoder (EcoVAE), to predict species distributions for the first time using large-scale, unstructured, and sparse occurrence data. To demonstrate the effectiveness of our framework, we trained a customized EcoVAE on a massive global dataset including

3

81    nearly 34 million georeferenced vouchered occurrences from the GBIF platform spanning

82    13,125 plant genera and 127,281 species. EcoVAE learns the patterns of global plant

83    distributions without any dependency on environmental variables. It efficiently processes

84    this large-scale dataset with fast computational times, accommodates varying rates of

85    missing data and biases, and enables the study of complex interactions within

86    communities. Remarkably, EcoVAE can accurately reconstruct full plant distributions

87    across all genera using as little as 20% of randomly selected occurrence records. We

88    further demonstrate the broad applicability of EcoVAE by applying it to 68 million

89    occurrence records of butterflies and 22 million records of mammals at both the genus and

90    species levels. Additionally, our model predictions offer an unsupervised approach to

91    assess collection completeness across different taxa on a global scale. These results

92    demonstrate the unprecedented capacity of deep learning methods to decode and predict

93    biodiversity patterns at global scales.

94    EcoVAE applies a unique masked approach to model the global species distributions using

95    well-curated occurrences data based on vouchered specimen records at the rank of genus

96    or species (Fig. 1a). We treat genera as the unit of inference in the modeling process due to

97    their computational efficiency and biological relevance. Genera represent coherent,

98    morphologically similar, and often monophyletic groups of species, providing a practical

99    compromise between taxonomic detail and manageable computational demands. The input

100   data was grouped into grids of 0.1', where plant observations within each grid were

101   summarized into vectors. The richness of genera per grid varied widely, ranging from $10^1$

102   to $10^3$ (Supplementary Fig. 1a). Our model consists of an encoder that learns a low-

103   dimensional representation of the input data and a decoder that reconstructs the presence

104   of genera per grid. For model training, we randomly masked 50% of the genera presence

105   data and the model was trained to predict these masked genera based on the remaining

106   observed data (Fig. 1b). This process allows the model to interpolate sparse observations

107   and estimate the true, unobserved plant distributions.

108   We evaluated the performance of our model in three randomly selected regions in North

109   America, Europe, and Asia, and applied the remaining global data for training (Fig. 1c,

110   Supplementary Table 1). Our model demonstrated a 10-fold increase in computational

4

111   speed compared to popular traditional SDMs, including logistic regression and random

112   forest, when predicting the distribution of a single genus using all other genera as input

113   (Fig. 1d). This difference is even greater when predicting distributions of multiple genera

114   (Fig. 1e).

115   We further assess the model's accuracy in prediction. Here, we calculated the predicted

116   genera counts per grid for the test regions and compared them to the actual observations.

117   The results demonstrated very high Pearson correlation coefficients of 0.98, 0.99 and 0.99

118   for test regions in North America, Europe, and Asia, respectively. At the species level, our

119   model achieves correlation coefficients of 0.95, 0.98, and 0.98 across the three regions

120   (Supplementary Fig. 2). We also calculated the total number of genera present in each grid

121   and observed high correlations between the model's predictions and the actual data (Fig.

122   1f).

123   Next, we used the Area Under the Receiver Operating Characteristic (AUROC) curve to

124   evaluate EcoVAE's accuracy in modeling the distribution of each genus. For the masked

125   genera, the mean AUROC was 0.82 for North America, 0.83 for Europe, 0.85 for Asia, which

126   demonstrates the robust performance of our model to infer missing information from

127   incomplete datasets (Fig. 1g). For example, *Lonicera* has a localized distribution in North

128   America and our model correctly predicted this pattern despite this genus being masked in

129   the input, with an overlap rate of 0.90 (Methods). Similar performance was observed for

130   *Lamium* in Europe (0.89) and for *Rhus* in Asia (0.80), which exhibit more scattered

131   distribution patterns (Fig. 1i). It is important to note that the three regions we selected

132   randomly differ substantially in plant distributions, area size, and genera counts per grid

133   (Supplementary Fig. 1b-d, Supplementary Table 1). Nevertheless, our model performed

134   equally well across them, which highlights the wide applicability of EcoVAE to diverse

135   geographic contexts. Herbarium specimen records represent a sparse sampling of actual

136   plant distributions, and data completeness varies significantly across regions [28,29]. To

137   address this inherent limitation, we analyzed the impact of data sparsity on our model's

138   performance. We tested our model using only 1%, 5%, 10%, 20%, and 30% of the input

139   genera, and evaluated its performance based on the AUROC for the remaining genera. With

140   only 1% of the input data, the model's performance was relatively low and the mean

141    AUROC was 0.56. However, increasing the input to only 5% improved the mean AUROC to

142    0.68. The mean AUROC further rose to 0.77 when we used 20% of the input data, close to

143    the performance seen with 50% (AUROC of 0.81). These results demonstrate that our

144    trained model can effectively use as little as 20% of the available data to reconstruct the full

145    generic distribution with high precision (Fig. 1h).

146    We extended our modeling framework to other major clades with high conservation values,

147    i.e., butterflies and mammals, and evaluated its performance using a similar approach with

148    three test regions (Supplementary Table 2). For butterflies, we found that our model

149    achieves high accuracy in predicting the number of genera per grid, with the Pearson's

150    correlation coefficients of 0.96, 0.99, 0.80 for genera counts per grid for the test regions

151    (Supplementary Fig. 3). The AUROC scores for genus-level predictions were 0.79, 0.84, and

152    0.75 for these regions. At the species level, the model achieved comparable results for

153    North America and Europe, but the AUROC decreased to 0.68 for Asia, which may reflect its

154    incompleteness of vouchered occurrences at the species level. For mammals, the model

155    performed best in the test region of North America at both genus and species levels

156    (Supplementary Fig. 4). In contrast, the sparser data in Asia posed challenges for

157    reconstructing full species distributions. Overall, our results demonstrate that EcoVAE

158    generalizes effectively across diverse taxa and geographies.

159    One important application of species modeling is interpolating occurrences where data

160    were lacking. We assume that the prediction error of EcoVAE reflects the completeness of

161    the occurrence records: if the records are incomplete, the model will struggle to

162    reconstruct the input data effectively. We estimated the prediction error globally

163    (Methods) and found that regions with high prediction error overlap with known

164    "darkspots" of biodiversity collection [30,31]. For example, the highest prediction errors for

165    plants were observed in South Asia, Southeast Asia, the Middle East, and Central Africa.

166    South America showed higher prediction errors compared to North America (Fig. 2b).

167    Notably, despite generally sparse records from high-latitude regions, the prediction error

168    remained low, suggesting that the occurrence records in these areas are nearly complete,

169    which allows the model to reflect true species distributions (Fig. 2b) more accurately. For

170    butterflies, the highest prediction errors were observed in South America and parts of

171 Southeast Asia (Supplementary Fig. 5). Interestingly, Middle Africa exhibited low

172 prediction errors, in contrast to the patterns observed for plants. For mammals, the

173 prediction error was generally smaller, likely due to the low genus diversity. However,

174 regions in South America and Central Asia displayed comparatively high prediction errors,

175 highlighting the need for further investigation efforts in these regions.

176 We then assessed the interpolation power of EcoVAE on i.) a region in southeastern North

177 America with relatively incomplete herbarium records but rich observation data from

178 iNaturalist, and ii.) a region in South Asia with sparse online occurrence records of both

179 kinds (Fig. 2a).  We applied the same model structure to train a full global model based on

180 all available plant voucher records and generated the new predictions in this test region

181 (Methods). In North America, we calculated an overlap index for each genus with

182 iNaturalist observations, defined as the ratio of predictions that are absent from input data

183 but present in the iNaturalist data. We found that our model performed best for genera

184 with a moderate number of observations, while abundant data results in diminishing

185 returns from interpolation (Supplementary Fig. 6). Using the genus *Sassafras* as an

186 example, we found that the new predictions largely overlap with the iNaturalist data (Fig.

187 2b).

188 For regions like South Asia, both georeferenced vouchered and observational datasets are

189 sparser. We selected genera that showed significant expansion in our model's new

190 predictions. For example, the *Desmodium* genus in the Fabaceae family only have

191 vouchered specimens in the eastern Himalayan/Nepal region in GBIF, but our model

192 predicts its much wider distribution across western and southern India (Fig. 2d), which

193 aligns better with field surveys and floristic investigations [32,33]. Similarly, for *Melicope* in

194 Rutaceae, the original observations were localized in southern India, but our new

195 predictions included occurrences in Myanmar and lowlands of Nepal, which is also

196 confirmed by third-party observations (Fig. 2d) [34,35]. For *Adonis*, our new predictions

197 suggested a broad distribution across the Himalayan region, consistent with various local

198 floras and checklists describing the widespread nature of the genus from Pakistan to

199 temperate regions in China (Fig. 2d) [36]. These results highlight the power of our model to

200 uncover plant distribution patterns in regions where observational data are limited.

7

201   Another important aspect of distribution modeling is understanding the community

202   response to changing distributions of organisms within it. Here, we interrogated our full

203   global model to study genus-to-genus interactions. We selected a test region in Australia

204   with abundant occurrences data and a good representation of major biomes

205   (Supplementary Fig. 7, Supplementary Table 1). We conducted *in silico* perturbation

206   analysis, where each grid cell in the region was artificially altered by introducing a genus to

207   areas where it was previously absent. By comparing the perturbed predictions to

208   unperturbed models (Fig. 2e), we assessed the invasive potential of one genus on others'

209   distributions. We focused on statistically significant interactions for downstream analysis

210   (Methods, Supplementary Fig. 8).

211   Examination of our genus network revealed that genera with high out-degree (those

212   influencing others significantly) tend to have low in-degree (being influenced by others),

213   suggesting asymmetric interactions [37] (Fig. 2f). Genera with broader ranges tend to interact

214   with a larger number of genera (Supplementary Fig. 9). We revealed that certain families,

215   including Poaceae, Cyperaceae, and Amaranthaceae, are more sensitive to disturbance in

216   the study region (Fig. 2g). These families were significantly influenced by members from

217   Poaceae, Asteraceae, and Fabaceae, which are globally well-represented in naturalized and

218   invasive floras (Supplementary Fig. 10) [38]. Thus, our model reveals patterns of genus

219   interactions, providing insights into community dynamics that may not be directly

220   observable in the original co-occurrence records.

221   In this work, we present EcoVAE, a generative deep learning framework for modeling

222   global plant distributions with high precision and speed. EcoVAE evaluates global sampling

223   efforts, interpolates distributions, and reveals interspecies interactions via *in silico*

224   perturbation analyses, offering novel applications for biodiversity exploration and

225   monitoring. It demonstrates that species distributions can be reconstructed using co-

226   occurrence information alone, even with incomplete data, capturing ecological patterns

227   often missed by traditional approaches. EcoVAE complements current SDMs by providing a

228   scalable framework for global analyses that can guide more targeted ecological studies. For

229   instance, it can identify under-sampled regions or unexpected patterns, directing SDM

230   efforts and field surveys to areas most in need of investigation. Furthermore, while we have

231 demonstrated the high performance of EcoVAE on taxa such as plants, butterflies, and

232 mammals, it can be easily extended to other taxa including birds and invertebrates. We

233 envision that EcoVAE will advance biodiversity investigations, especially in under-sampled

234 regions with limited environmental data, and ultimately support global biodiversity

235 monitoring efforts aligned with the Convention on Biological Diversity [39]. Future

236 integration of additional data, such as geographic or climate variables, could potentially

237 improve performance and reveal insights into organism distributions and environmental

238 change.

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

## Methods

### Data preprocessing

We downloaded the world plant, butterfly and mammal distribution data from the GBIF platform (Supplementary Note 1). We used the R package "CoordinateCleaner" [40] to remove records located in the sea, on country or major area centroids, capitals, or in major biodiversity facilities. We modeled the distribution both at the genus and species level. The cleaned observation dataset includes 33.8 million observations of plant, 67.6 million observations of butterfly and 21.9 million observations for mammal (Supplementary Table 2). The grid size was set to 0.1' x 0.1' and we summed all the observation data in each grid. For any genus or species with more than one observation, we set the value to 1 (existence) in contrast to 0 (non-existence). For plant observations, we only kept genera or species that occurs in over 20 grid cells and grid cells with more than 5 different genera. The finalized data contains 277,133 grids, while the three test regions include 11,292, 3,878 and 1,210 grids for North America, Europe, and Asia respectively. For butterfly and mammal observations, we only kept genera or species that occur in more than 5 grid cells, and grid cells with more than 1 genus or species.

### Model structure

We developed an Ecological Variational Autoencoder (EcoVAE) model which aims to reconstruct the full plant distribution based on partial observations. The core VAE architecture comprises an encoder and a decoder.

1. Encoder function: the encoder is implemented as a sequential network with two hidden layers of Gelu activated linear transformations (dimension: 128). The encoder maps the input data (e.g., dimension: 13,125 for plant genera) into a latent space characterized by mean ($\mu$) and log variance ($log\ (\sigma^2)$) parameters (dimension: 32).

$$(\mu, log\ log\ (\sigma^2)\ ) = f_{encoder}(x) = GELU\big(W_2\big(GELU(W_1 x + b_1)\big) + b_2\big)$$

10

282      Where W1, b1, W2, b2 are the weights and biases of the two hidden layers, and

283      GELU denotes the Gaussian Error Linear Unit activation function.

284    2.  Reparameterization: the latent representation is obtained through a

285      reparameterization step that ensures differentiability by sampling from a Gaussian

286      distribution parameterized by these $\mu$ and $log\ (\sigma^2)$.

287
$$z = \mu + \sigma \odot \epsilon$$

288      where $\epsilon$ follows normal distribution, $\epsilon \sim N(0,\ I)$.

289    3.  Decoder function: the decoder, mirroring the encoder's structure, reconstructs the

290      input data from the latent space, aiming to minimize reconstruction error. The

291      output dimension is equivalent to the input dimension (13,125), ensuring that the

292      reconstructed output mirrors the input feature set.

293
$$\hat{x} = f_{decoder}(z) = GELU\big(W_4\big(GELU(W_3 z + b_3)\big) + b_4\big)$$

294      Where W3, b3, W4, b4 are the corresponding weights and biases for the decoder

295      layers.

296  This mathematical framework enables the EcoVAE to compress high-dimensional data into

297  a lower-dimensional latent space and subsequently reconstruct the original data with

298  minimized reconstruction error.

299

300  **Model training and evaluation**

301  During the model training process, we used a unique masking strategy where 50% of input

302  data are randomly set to zero to simulate missing data scenarios. The model uses a

303  weighted reconstruction loss function that is based on mean squared error (MSE). A

304  weighting factor of 0.5 modifies the contribution of masked and unmasked genera to the

305  reconstruction loss, providing a balanced approach to learning from both visible and

306  obscured portions of the data. We used the Adam optimizer with a fixed learning rate of

307  0.001, and the models were trained for 15 epochs with a batch size of 512.

308    To evaluate the model's performance, we excluded data from three randomly selected

309    regions in Europe, North America, and China, during the training phase (Supplementary

310    Table 1). After model training, we conducted a masking procedure where we randomly

311    selected 50% of all genera or species from these regions and set their corresponding

312    observational data to zero. These masked data were then used as inputs to evaluate the

313    model's ability to reconstruct the observational data for the masked genera or species from

314    the unmasked genera. We quantitatively measured the model's performance using the Area

315    Under the Receiver Operating Characteristic (AUROC) and MSE metrics (Supplementary

316    Fig. 2). For each genus or species in the test regions, we selected an equal number of top

317    predicted grids as in the original data. We calculated the overlap rate as the fraction of

318    predictions that had a true occurrence record either within that grid or in neighboring

319    grids.

320

**Model benchmarking**

322    To evaluate the time efficiency of our EcoVAE model in processing large-scale input data,

323    we conducted a series of benchmarking experiments in  comparison with two popular SDM

324    methods, i.e., random forest and logistic regression [18]. We tested the model's performance

325    under varying input dimensions by randomly selecting genera from the input genus

326    presence matrix. For these experiments, we used a masking strategy where 80% of the

327    input columns were masked to simulate missing data. A binary mask was generated using a

328    probability threshold proportional to the desired masking percentage, ensuring that

329    columns in the input data were set to zero with the defined probability. We used the

330    training and test data split as previously described. A random genus with presence in at

331    least five grids was randomly chosen as the prediction target: the model should predict the

332    presence of this genus across all grids in the training and test data based on the masked

333    input matrix. We trained the random forest classifier using the

334    *sklearn.ensemble.RandomForestClassifier* function with default parameters. We trained the

335    logistic regression model with the *sklearn.linear_model. LogisticRegression* function with

336    the following parameters: max_iter=1000, solver='saga', penalty='l1'. Time measurements

337    were recorded from the start of data preparation to the completion of the prediction phase.

338    The time consumption was calculated over 10 iterations to benchmark the models' time

339    efficiency.

340

341    **Model application**

342    **-Data interpolation**

343    We trained the EcoVAE model as previously described on all available occurrence data for

344    interpolating unobserved plant distribution (full global model). To evaluate the model's

345    prediction error, we used the unmasked global data as input and applied a threshold to

346    binarize the output genus presence matrix, ensuring that the total number of occurrences

347    was doubled. For each grid, we then calculated the ratio of observed genera not

348    represented in the output matrix, which we defined as the prediction error. Based on the

349    global distribution of prediction error, we selected two regions to evaluate the

350    performance of data interpolation, i.e., North America and South Asia. For these regions, we

351    used a similar strategy to generate the binarized output for downstream analysis.

352    For the North America region, we used observation data collected from iNaturalist to verify

353    the accuracy of model prediction instead of traditional data splitting method. The ratio

354    between iNaturalist observational data and vouchered specimen data is 7:1 and we would

355    expect that iNaturalist data have a better geographic coverage than herbarium specimens

356    for many species. We calculated overlapping rate between predicted species occurrences

357    and actual observations to quantify the performance of our prediction. For each species, we

358    first extracted the observed and predicted occurrences based on the genus index. For both

359    datasets, we retained only the presence points. The observed points were buffered by 0.1

360    degrees to account for geographic uncertainties. We then converted the presence data into

361    spatial objects using the "sf" package in R [41]. Overlap between predicted occurrences and

362    observed occurrences, as well as overlap between predicted and original input

363    occurrences, was calculated using spatial intersections (`st_intersects`). The key metric, the

364    overlapping rate, was calculated by dividing the number of predicted points that

365    overlapped with observed points but not with original input points by the total number of

13

366    predicted points. This rate reflects the proportion of new predicted occurrences that align

367    with observed data but were not part of the input data, providing a measure of prediction

368    accuracy and novelty. We also compared the predicted overlapping rate with the

369    overlapping rate calculated between the same number of randomly generated points in the

370    study area and observations.

371    For the South Asia region, we assessed the occurrence of each genus both before and after

372    data interpolation. We focused on genera that initially occurred in more than 5 grids and

373    whose distribution region has expanded most for downstream analysis. Due to the lack of

374    georeferenced observational data in this region, we compared our prediction with the

375    distribution described in plant atlas and related literature.

376

377    **-Simulation of genera interaction**

378    To simulate the impact of a specific genus i on all other genera within a targeted region, we

379    initially identified all grid cells lacking genus i. The observational data from these grids

380    were utilized as input for the model, and the corresponding reconstructed data served as

381    the background dataset (x_background). Then we introduced observations of genus i into

382    these grids and generated perturbed model outputs (x_perturb). By comparing the plant

383    distributions between x_background and x_perturb, we were able to identify genera that

384    exhibited significant changes, thereby quantifying the ecological influence of genus i on the

385    plant community dynamics within the region.

386    To assess species interactions after species additions, we first fit linear regression models

387    to compare grid numbers of all genera before and after addition of a specific genus i.

388    Specifically, we used the 'lm' function in R to model the relationship between grid numbers

389    before and after addition of genus i. For each model, we calculated 99.99% confidence

390    intervals using the 'predict' function with the interval parameter of "prediction" and level

391    parameter of 0.9999. We defined the significant interaction between genus i and j if the

392    predicted grid number for genus j falls outside the bounds of the confidence intervals after

393    addition of genus i. In such circumstances, we identified j as an "outlier" and defined it as a

394    "sensitive genus". Z-scores were calculated for each genus by normalizing the residuals,

14

395 computed as the difference between actual and predicted values. We then performed

396 frequency analysis of the impactful and vulnerable genera based on all significant

397 interactions. To further explore genus interactions at the family level, we utilized the

398 "plantlist" package [42] to classify genera into families and analyzed the proportion of

399 sensitive genera within each family. We selected the most sensitive family based on the

400 following criteria: it includes more than 5 genera and at least 35% of the genera are

401 classified as "sensitive" (significantly impacted by at least one other genus).

402

## Author Contributions

404 YY and SB conceptualized the study, conceived EcoVAE, collected and analyzed the data. YY

405 and SB wrote the manuscript with key contributions from CCD. All authors approved the

406 manuscript.

407

## Acknowledgements

412

## References

414 1. Guisan, A. *et al.* Predicting species distributions for conservation decisions. *Ecology*

415    *Letters* **16**, 1424–1435 (2013).

416 2. Franklin, J. Species distribution models in conservation biogeography: developments and

417    challenges. *Diversity and Distributions* **19**, 1217–1223 (2013).

418 3. Franklin, J. Species distribution modelling supports the study of past, present and future

419    biogeographies. *Journal of Biogeography* **50**, 1533–1545 (2023).

420 4. Elith*, J. *et al.* Novel methods improve prediction of species' distributions from

15

421  occurrence data. *Ecography* **29**, 129–151 (2006).

422 5. Chollet Ramampiandra, E., Scheidegger, A., Wydler, J. & Schuwirth, N. A comparison of

423  machine learning and statistical species distribution models: Quantifying overfitting

424  supports model interpretation. *Ecological Modelling* **481**, 110353 (2023).

425 6. Franklin, J. Species distribution modeling. in *Mapping Species Distributions: spatial*

426  *inference and prediction* 3–20 (Cambridge University Press, Cambridge, United Kingdom

427  and New York, NY, USA, 2009).

428 7. Zurell, D. *et al.* Benchmarking novel approaches for modelling species range dynamics.

429  *Global Change Biology* **22**, 2651–2664 (2016).

430 8. Franklin, J., Serra-Diaz, J. M., Syphard, A. D. & Regan, H. M. Big data for forecasting the

431  impacts of global change on plant communities. *Global Ecology and Biogeography* (2016)

432  doi:10.1111/geb.12501.

433 9. Runting, R. K., Phinn, S., Xie, Z., Venter, O. & Watson, J. E. M. Opportunities for big data in

434  conservation and sustainability. *Nat Commun* **11**, 2003 (2020).

435 10. Di Cecco, G. J. *et al.* Observing the Observers: How Participants Contribute Data to

436  iNaturalist and Implications for Biodiversity Science. *BioScience* **71**, 1179–1188 (2021).

437 11. Davis, C. C. The herbarium of the future. *Trends in Ecology & Evolution* **38**, 412–423

438  (2023).

439 12. Hedrick, B. P. *et al.* Digitization and the Future of Natural History Collections.

440  *BioScience* **70**, 243–251 (2020).

441 13. Araújo, M. B. & Luoto, M. The importance of biotic interactions for modelling species

442  distributions under climate change. *Global Ecology and Biogeography* **16**, 743–753

443  (2007).

444 14. Van der Putten, W. H., Macel, M. & Visser, M. E. Predicting species distribution and

445  abundance responses to climate change: why it is essential to include biotic interactions

446  across trophic levels. *Philosophical transactions of the Royal Society of London. Series B,*

447  *Biological sciences* **365**, 2025–2034 (2010).

448 15. Thuiller, W. *et al.* Navigating the integration of biotic interactions in biogeography.

449  *Journal of Biogeography* **51**, 550–559 (2024).

450 16. Wisz, M. S. *et al.* The role of biotic interactions in shaping distributions and realised

451  assemblages of species: Implications for species distribution modelling. *Biological*

452    *Reviews* **88**, 15–30 (2013).

17.    Guisan, A. *et al.* Sensitivity of predictive species distribution models to change in grain size. *Diversity and Distributions* **13**, 332–340 (2007).

18.    Valavi, R., Guillera-Arroita, G., Lahoz-Monfort, J. J. & Elith, J. Predictive performance of presence-only species distribution models: a benchmark study with reproducible code. *Ecological Monographs* **92**, e01486 (2022).

19.    Beck, J., Böller, M., Erhardt, A. & Schwanghart, W. Spatial bias in the GBIF database and its effect on modeling species' geographic distributions. *Ecological Informatics* **19**, 10–15 (2014).

20.    Garcia-Rosello, E., Gonzalez-Dacosta, J., Guisande, C. & Lobo, J. M. GBIF falls short of providing a representative picture of the global distribution of insects. *Systematic Entomology* **48**, 489–497 (2023).

21.    Pender, J. E. *et al.* How sensitive are climatic niche inferences to distribution data sampling? A comparison of Biota of North America Program (BONAP) and Global Biodiversity Information Facility (GBIF) datasets. *Ecological Informatics* **54**, 100991 (2019).

22.    Dormann, C. F. *et al.* Collinearity: A review of methods to deal with it and a simulation study evaluating their performance. *Ecography* **36**, 027–046 (2013).

23.    Ouyang, L. *et al.* Training language models to follow instructions with human feedback. Preprint at https://doi.org/10.48550/arXiv.2203.02155 (2022).

24.    Ramesh, A. *et al.* Zero-Shot Text-to-Image Generation. Preprint at https://doi.org/10.48550/arXiv.2102.12092 (2021).

25.    Weaver, W. N., Ruhfel, B. R., Lough, K. J. & Smith, S. A. Herbarium specimen label transcription reimagined with large language models: Capabilities, productivity, and risks. *American J of Botany* **110**, e16256 (2023).

26.    Lopez, R., Regier, J., Cole, M. B., Jordan, M. I. & Yosef, N. Deep generative modeling for single-cell transcriptomics. *Nat Methods* **15**, 1053–1058 (2018).

27.    Kingma, D. P. & Welling, M. Auto-Encoding Variational Bayes. Preprint at https://doi.org/10.48550/arXiv.1312.6114 (2022).

28.    Daru, B. H. *et al.* Widespread sampling biases in herbaria revealed from large-scale digitization. *New Phytol* **217**, 939–955 (2018).

483    29.    Daru, B. H. & Rodriguez, J. Mass production of unvouchered records fails to

484           represent global biodiversity patterns. *Nat Ecol Evol* **7**, 816–831 (2023).

485    30.    Daru, B. H. Predicting undetected native vascular plant diversity at a global scale.

486           *Proc. Natl. Acad. Sci. U.S.A.* **121**, e2319989121 (2024).

487    31.    Ondo, I. *et al.* Plant diversity darkspots for global collection priorities. *New*

488           *Phytologist* **244**, 719–733 (2024).

489    32.    Joshi, B. R., Hakim, M. M. & Patel, I. C. The biological active compounds and biological

490           activities of Desmodium species from Indian region: a review. *Beni-Suef Univ J Basic Appl*

491           *Sci* **12**, 1 (2023).

492    33.    Desmodium gangeticum (L.)DC. Species. *India Biodiversity Portal*

493           https://indiabiodiversity.org/species/show/229507.

494    34.    Wood, K. R., Appelhans, M. S. & Wagner, W. L. Melicope oppenheimeri, section Pelea

495           (Rutaceae), a new species from West Maui, Hawaiian Islands: with notes on its ecology,

496           conservation, and phylogenetic placement. *PK* **69**, 51–64 (2016).

497    35.    Hartley TG. On the Taxonomy and Biogeography of Euodia and *Melicope* (Rutaceae).

498           in *Allertonia* vol. 8(1) (National Tropical Botanical Garden, Lawaʻi, Kauaʻi, Hawaiʻi.,

499           2001).

500    36.    Fu, D. Adonis, in *Flora of China* **6**: 389-391. (2001).

501    37.    Bascompte, J., Jordano, P. & Olesen, J. M. Asymmetric Coevolutionary Networks

502           Facilitate Biodiversity Maintenance. *Science* **312**, 431–433 (2006).

503    38.    Pyšek, P. *et al.* Naturalized alien flora of the world: species diversity, taxonomic and

504           phylogenetic patterns, geographic distribution and global hotspots of plant invasion.

505           *Preslia* **89**, 203–274 (2017).

506    39.    Díaz, S. *et al.* Set ambitious goals for biodiversity and sustainability. *Science* **370**,

507           411–413 (2020).

508    40.    Zizka, A. *et al.* CoordinateCleaner : Standardized cleaning of occurrence records from

509           biological collection databases. *Methods in Ecology and Evolution* **10**, 744–751 (2019).

510    41.    Pebesma, E. Simple Features for R: Standardized Support for Spatial Vector Data.

511           *The R Journal* **10**, 439 (2018).

512    42.    Zhang, J. Plantlist: Looking Up the Status of Plant Scientific Names Based on The

513           Plant List Database. (2019). R package version

514    (4.3.3) https://github.com/helixcn/plantlist/

515

**Code Availability**

517    Our trained model and codes are available from GitHub:

518    https://github.com/lingxusb/EcoVAE

519

**Additional Information**

521    Supplementary information

522

523

524

525

526

527

528

529

530

531

532

533

534 **Figure Legends**

535 **Figure 1. Framework and evaluation of EcoVAE model performance. a,** Schematic

536 representation of the model training and application pipeline. **b,** Overview of the model

537 evaluation process. **c,** Map showing the locations of the three testing regions. **d,**

538 Comparison of time consumption between two machine learning methods and EcoVAE

539 with the increase of input dimension. **e,** Comparison of time consumption between two

540 machine learning methods and EcoVAE with the increase of output dimension. **f,**

541 Correlation between observed genera counts per grid (or observed grid counts per genus,

542 upper panels) and predicted genera counts per grid (or predicted grid counts per genus,

543 lower panels) across the three testing regions, with high Pearson correlation values. The

544 black dashed lines indicate identity lines. **g,** Density plot showing the Area Under the

545 Receiver Operating Characteristic (AUROC) for three testing regions. **h,** Relationship

546 between the ratio of masked data and AUROC values. **i,** Comparison of the distribution of

547 observations (upper panels) and predictions (lower panels) for randomly selected genus

548 within each test region.

549

550 **Figure 2. Applications of EcoVAE model. a,** Schematic illustration of the interpolation

551 process using EcoVAE. **b,** Global distribution of relative collection completeness,

552 represented by the value of prediction error of EcoVAE. Darker color represents lower

553 prediction error and higher completeness, while lighter color represents higher prediction

554 error and lower completeness. **c,** Comparison between original herbarium specimen

555 records and EcoVAE interpolation results of genus *Sassafras* in North America. **d,**

556 Interpolation results of three example genera in South Asia. Gray dots show the

557 distribution of all georeferenced vouchered occurrences in the study area. **e,** Schematic

558 illustration of studying community interactions using EcoVAE model. **f,** Relationship

559 between the number of outdegree and indegree for the genus-to-genus interactions. Each

560 dot represents a single genus. **g,** Log number of sensitive genera across the most sensitive

561 plant families identified by EcoVAE.

562

563