

Spatial occurrence records and distributions of tropical Asian butterflies

*Eugene Yu Hin Yau¹, *Emily E. Jones¹, Toby Pak Nok Tsang^{1,2}, Shuang Xing^{1,3}, Richard T. Corlett⁴, Patrick Roehrdanz⁵, David J. Lohman⁶, Adam Kai Chi Lee¹, Catherine Wai Ching Hai¹, Shawan Chowdhury^{7,8,9,10}, Jane K. Hill¹¹, Jade A. T. Badon¹², Cheong Weei Gan¹³, Yves Basset¹⁴, I-Ching Chen¹⁵, Suzan Benedick¹⁶, Anuj Jain^{13,17}, Tiffany L.T. Ki^{11,18}, Krushnamegh Kunte¹⁹, Akihiro Nakamura²⁰, Lien Van Vu²¹, Sarah A. Scriven¹¹, Alice C. Hughes¹, Timothy C. Bonebrake^{1#}

*authors contributed equally, #corresponding author

¹School of Biological Sciences, The University of Hong Kong, Pokfulam, Hong Kong SAR, CN

²The University of Toronto Scarborough, 1265 Military Trail, Scarborough, ON M1C 1A4, CA

³School of Ecology, Shenzhen Campus of Sun Yat-sen University; Shenzhen 518107, China.

⁴Center for Integrative Conservation and Yunnan Key Laboratory for the Conservation of Tropical Rainforests and Asian Elephants, Xishuangbanna Tropical Botanical Garden, Chinese Academy of Sciences; Honorary Research Associate, Royal Botanic Gardens, Kew.

⁵Betty and Gordon Moore Center for Science, Conservation International, Arlington, VA, USA

⁶City College of New York, City University of New York, 160 Convent Avenue New York, NY 10031, USA; PhD Program in Biology, City University of New York, 365 Fifth Avenue, New York, NY 10016; Zoology Division, National Museum of Natural History, Rizal Park, Manila 1000, Philippines

⁷Institute of Biodiversity, Friedrich Schiller University Jena, Dornburger Straße 159, 07743 Jena, Germany

⁸Department of Biodiversity and People, Helmholtz Centre for Environmental Research – UFZ, Permoserstraße 15, 04318 Leipzig, Germany

⁹German Centre for Integrative Biodiversity Research (iDiv) Halle-Jena-Leipzig, Puschstraße 4, 04103 Leipzig, Germany

¹⁰Czech University of Life Sciences Prague, Faculty of Environmental Sciences, Prague, Czech Republic

¹¹Leverhulme Centre for Anthropocene Biodiversity, Department of Biology, University of York, York, YO10 5DD, UK

¹²Animal Biology Division, Institute of Biological Sciences, University of the Philippines Los Baños, Laguna 4031, Philippines

¹³Nature Society Singapore, 510 Geylang Road, 389466, Singapore

¹⁴Smithsonian Tropical Research Institute, Apartado 0843-03092, Balboa, Ancon, Panama

¹⁵Department of Life Sciences, National Cheng Kung University, Taiwan

¹⁶Faculty of Sustainable Agriculture, Universiti Malaysia Sabah, Locked Bag No. 3, 90509, Sandakan, Sabah, Malaysia

¹⁷bioSEA Pte Ltd., 68 Chestnut Avenue, 679521, Singapore

¹⁸Science Department, Natural History Museum, London SW7 5BD, United Kingdom

¹⁹National Centre for Biological Sciences (NCBS), Tata Institute of Fundamental Research (TIFR), GKVK Campus, Bellary Road, Bengaluru 560065, India

²⁰CAS Key Laboratory of Tropical Forest Ecology, Xishuangbanna Tropical Botanical Garden, Chinese Academy of Sciences, Menglun, China

50 ²¹Vietnam National Museum of Nature, Vietnam Academy of Science and Technology, 18
51 Hoang Quoc Viet, Cau Giay, Ha Noi, Vietnam

52
53 corresponding author(s): Timothy C. Bonebrake (tbone@hku.hk)

54
55

56 **Abstract**

57

58 Insect biogeography is poorly documented globally, particularly in the tropics. Recent
59 intensive research in tropical Asia, combined with increasingly available records from citizen
60 science, provides an opportunity to map the distributions of tropical Asian butterflies. We
61 compiled a dataset of 724,247 occurrences of 3,591 tropical Asian butterfly species by
62 aggregating records from GBIF (651,285 records), published literature (21,271), published
63 databases (37,695), and unpublished data (13,993). Here, we present this dataset and single-
64 species distribution maps of 1,520 species. Using these maps, along with records of the 2,071
65 remaining species, we identified areas of limited sampling (e.g., the Philippines, Myanmar,
66 and New Guinea) and predicted areas of high diversity (Peninsular Malaysia and Borneo).
67 This dataset can be leveraged for a range of studies on Asian and tropical butterflies,
68 including 1) species biogeography, 2) sampling prioritization to fill gaps, 3) biodiversity
69 hotspot mapping, and 4) conservation evaluation and planning. We encourage the continued
70 development of this dataset and the associated code as a tool for the conservation of tropical
71 Asian insects.

72

73

74 **Background & Summary**

75

76 Tropical Asia, home to multiple major global biodiversity hotspots, harbors a rich assemblage
77 of highly range-restricted endemic species¹. Unfortunately, reliable distribution data for many
78 species in this region are scarce². One prominent challenge for invertebrate conservation,
79 known as the Wallacean shortfall, stems from our inadequate knowledge of species
80 distributions³. Insufficient information on species distributions impedes the identification of
81 vulnerable species and the efficient allocation of conservation resources across regions and
82 species^{3,4}.

83

84 While recent global studies of butterfly biogeography have incorporated data from tropical
85 Asia^{5,6}, they have primarily relied on coarse, country-level data to examine biogeographic
86 patterns⁵⁻⁷. The distribution information summarized based on those data is largely influenced
87 by political boundaries rather than relevant ecological areas and is less ideal for identifying
88 important conservation/vulnerable areas, which requires fine-scale, biogeographic data with
89 low bias⁸. There have also been attempts to map spatial phylogenetic diversity using range
90 maps⁹, but the quality of such spatial analyses is highly dependent on the range maps used,
91 which often fail to capture distribution patterns at local scales, thereby limiting the resolution
92 of the spatial pattern of interest. Although fine-scale geographic distributions of several Asian
93 butterfly groups have been mapped (e.g., *Elymnias* in Wei et al.¹⁰; *Papilio* in Condamine et
94 al.¹¹; *Polyura* in Toussaint et al.¹²; range-restricted butterflies in Scriven et al.¹³), to date, no
95 unified, fine-scale distribution dataset has been produced for the entire region – despite the
96 importance of such a tool for examining patterns of diversity within this highly biodiverse
97 region^{1,6}. Existing locality data might not be readily accessible and frequently require
98 aggregation and standardization. Fine-grained information on species distributions is an
99 essential first step for understanding insect biodiversity patterns and conservation needs.

100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148

The creation of regional datasets of species distributions is aided by the recent development of large, open-source biodiversity data platforms such as the Global Biodiversity Information Facility (hereafter, GBIF), an online database that organizes crowd-sourced data from citizen science platforms, scientific literature, and specimen collections¹⁴. These data, however, often include large spatial biases due to uneven sampling and data mobilization efforts among regions^{14,15}. Even if available, much of the fine-scale biogeographic data that could be employed to reduce these biases remains buried in literature and regional databases, requiring concerted efforts to make it analysis-ready⁷. Without unified and standardized datasets, it is difficult to test macroecological and macroevolutionary questions¹⁶, produce high-quality species distribution models¹⁷, and identify effective conservation targets^{5,6,8,18}.

The process of mapping species distributions can be accomplished either through data-driven modeling or by relying on expert knowledge. Range maps (expert range maps) solely based on expert knowledge tend to overestimate active areas of species at the local scale^{15,17,19}. In addition, the quality of their source data, hence the uncertainty of the analysis, is often unknown¹⁶. The dependency of range maps on expert knowledge means this method is available to a small subset of well-studied species⁷. In contrast, data-driven distribution maps offer greater transparency and reproducibility^{18,20}. Modern modeling techniques allow the interpolation of potential distributions into areas for which primary data collection may not be possible, enabling the production of more detailed and reliable distribution maps^{3,21}. However, major data gaps exist for occurrence records of most taxa^{16,22}, particularly invertebrates, and the non-random distribution of these gaps necessitates careful treatment within models²³.

Species distribution maps facilitate the identification of species ranges and hotspots of diversity. This provides valuable insights for local conservation planning/prioritization^{24,25} and policy-making, paving the way for future investigations into butterfly biogeography⁵ and phylogeographic patterns²⁴. Specifically, species distribution maps can guide the allocation of conservation resources, inform the strategic design of protected areas in high-suitability/biologically diverse areas, and identify low-suitability areas in need of management^{25,26}, enabling effective conservation interventions. In conjunction with SDMs, occurrence datasets can help inform species reintroduction programs by identifying potentially suitable areas^{25,27} and optimal source populations²⁸, and expedite IUCN Red List assessment, which has poor species coverage in Asia. Additionally, applications of SDMs include the modeling of species and community-level responses to climate change^{24,27,29} and the assessment of extinction risks³⁰.

The need for species conservation is particularly acute in tropical Asia, defined broadly here to include South Asia and Southeast Asia (see Fig. 1). The area is home to over 20,000 islands, many of which were repeatedly connected and separated from adjacent landmasses during drastic sea-level fluctuations over the past 4 Mya³¹. This dynamic past led to the evolution of numerous species endemic to single islands or island groups, and as such this region hosts some of the world's greatest biodiversity – an estimated 15-25% of all well-studied terrestrial taxa and a large proportion of undescribed taxa^{32,33}. This highly biodiverse region is also one of the globe's most biologically threatened: it is estimated that 42% of Southeast Asia's biodiversity may be lost by 2100 as three quarters of its primary forests are lost to agriculture, urbanization, and mineral extraction^{32,34,35}.

149 We present a comprehensive dataset of tropical Asian butterflies, with more than half of the
150 records possessing high spatial accuracy (uncertainty < 10 km). This fills a major sampling
151 gap, given that Asia is poorly represented in global biodiversity data repositories^{15,22,36},
152 improved datasets are urgently needed to enable effective monitoring and management of
153 biodiversity across the region. Leveraging the data along with tailored species distribution
154 models (SDMs), we generate data-driven distribution maps at a resolution of 10 km x 10 km.
155 These maps enhance a fundamental understanding of butterfly macroecological patterns in
156 tropical Asia. Each butterfly species' distribution was individually modeled and, together
157 with buffered occurrence points of unmodelled species, employed to assess regional patterns
158 of species diversity. Combined with species distribution models, our aggregated data
159 advances knowledge of butterfly macroecology and facilitates evidence-based decision-
160 making for butterfly conservation in tropical Asia.

161

162

163 **Methods**

164

165 *Occurrence data*

166 We manually extracted GBIF records for tropical Asian Papilionoidea (Lepidoptera:
167 Nymphalidae, Papilionidae, Lycaenidae, Pieridae, Hesperidae, Riodinidae; -11.426 – 35.64
168 N, 67.588 – 174.990 E) for the years 1970-present on 15 April 2024 (Derived dataset
169 GBIF.org³⁷). The geographical extent of the study area was selected to encompass northern
170 temperate Asia to secure sufficient data to capture the full niche breadth of all species in the
171 subsequent SDMs. We included presence records derived from human observation, preserved
172 specimens, material samples, or literature, provided they had associated coordinates. We
173 omitted all records with >100,000 m coordinate uncertainty, so-called “fuzzy” taxon matches,
174 and records for which the scientific name was missing or incomplete unless nomenclature
175 could be extracted using a BOLD identifier (boldsystems.org/). This resulted in a final
176 number of GBIF records equalling 651,285.

177

178 Roughly 73% (472,714) of these records are ‘Research-grade’ observations from iNaturalist.
179 Information on how this designation is made is available at GBIF.org. The accuracy of
180 opportunistically collected data from crowd-sourced platforms like GBIF is often diminished
181 due to misidentifications, taxonomic, spatial, and temporal biases, as well as uneven
182 taxonomic validation due to lack of standardized reference data^{14,38-41}. Given these potential
183 issues, and to fill geographic gaps, we supplemented these GBIF data with expert data
184 (coauthor datasets, published literature) and harmonized binomials to a single expert dataset
185 (Lamas, 2015. Catalogue of the butterflies (Papilionoidea), available from the author.; see
186 below).

187

188 We extracted data from the B2D2 Database of Butterflies for Borneo provided by JKH/the
189 Darwin Initiative (n = 19,417), a dataset for Bangladesh provided by SC (Chowdhury et al.⁴²;
190 n = 18,278), and unpublished datasets from coauthors AN, DJL, LVV, TK, and YB (n =
191 13,993). For geographic regions with relatively few records (e.g., China, Myanmar, Thailand)
192 and for species with < 10 records, we conducted targeted searches of post-1970 published
193 literature on Google Scholar in English and Chinese (simplified and traditional) (genus OR
194 genus + species + country name), producing an additional 21,271 records. Although some
195 publications lacked collection dates for records (e.g., checklists), we assume that the
196 inclusion of species in recent publications is indicative of species' current localities.

197

198 For all records in published sources, we extracted coordinates, locality name, locality type
 199 (e.g., exact coordinates, city, national park, island, or province), country, and year of record
 200 (where available). If exact coordinates were not provided by the source, we used Google
 201 Earth Pro (v7.3.6.9345) to estimate the locality centroid for any record provided at the
 202 province level or below (e.g., national park or city). For records from islands ≤ 100 km at the
 203 widest dimension (e.g., localities within the Philippines and Indonesia), we estimated the
 204 island or archipelago centroid. If a range of coordinates was provided (e.g., records from The
 205 Butterflies of Vietnam), we selected a point within the range. Data sources for all records are
 206 provided in the reference column in Occurrence Records of Tropical Asian Butterflies: 1970-
 207 2024 (<https://doi.org/10.6084/m9.figshare.25037645>).

208

209 Final binomial harmonization, validation, and authority assignment were conducted by DJL
 210 using a taxonomic reference prepared by Gerardo Lamas (Lamas, 2015). Family names were
 211 aligned by hand to GBIF.

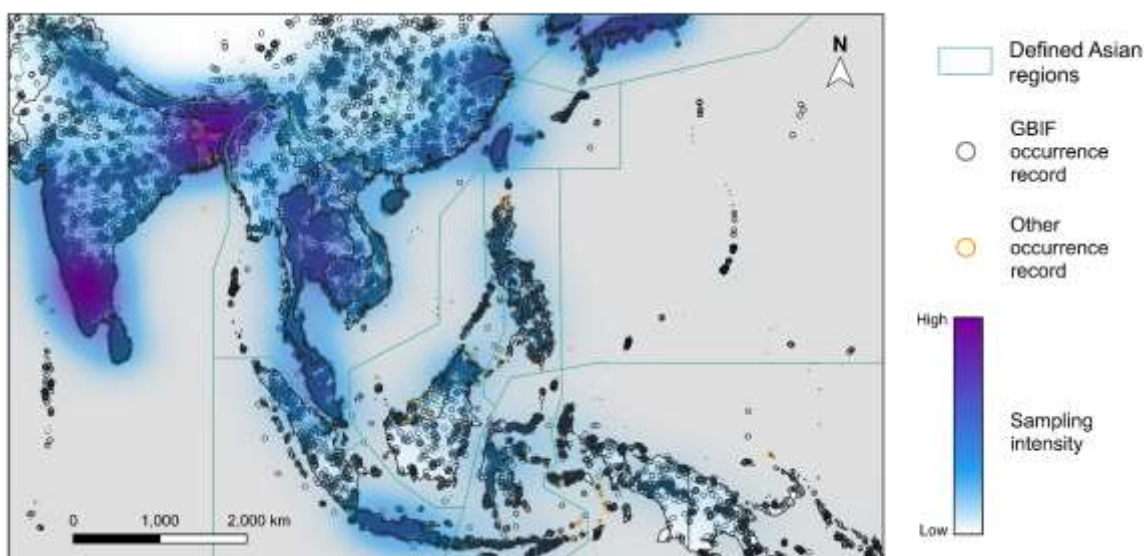
212

213 The resulting database consists of 724,247 occurrence records for 3,591 species from 546
 214 genera. These records represent approximately 20% of all described butterfly species
 215 globally^{43,44} (17,280-17,500 spp.; but see Pinkert et al., 2022⁵). Records of Nymphalidae
 216 (313,899; 1,324 spp.) comprise 43% of the dataset, followed by Lycaenidae (968 spp.;
 217 147,277 records), Papilionidae (264 spp.; 101,500 records), Pieridae (405 spp.; 97,120
 218 records), Hesperidae (611 spp.; 60,460 records), and Riodinidae (22 spp.; 3,991 records).

219

220 Of the 3,591 species in the database, 1,580 ($< 31\%$) are represented by ≥ 10 records within
 221 the extent of 36° N to 10° S and 69° E to 161.6° E that are >10 km apart (see details on
 222 distribution modeling below). Most occurrence records are concentrated in a limited number
 223 of regions, for example, India (28.34% of all data), Taiwan (13.75% of all data), Singapore
 224 (8.92% of all data), Hong Kong (8.30% of all data), and Malaysia (7.16% of all data) (Fig. 1).
 225 Equatorial regions together with southern China are relatively underrepresented in our
 226 dataset. As much of the data is derived from GBIF, which contains a large proportion of
 227 citizen science data, we observed a clustering of our data in areas of high population and a
 228 general lack of data in more inaccessible regions.

229



230

231

232 **Figure 1.** Distribution of GBIF and other occurrence records in our study area. Sampling
 233 intensity was estimated by running kernel density on the coordinates of all available
 234 occurrence data of every species. Regions of Asian landmasses based on the ecoregions and
 235 biogeographic realms as revised by Dinerstein et al.⁴⁵, as well as Wallace's Line, Huxley's
 236 Line, and Weber's Line.

237

238 *SDM methods and results*

239 Five algorithms, Generalized Linear Model (GLM), Maximum Entropy (MaxEnt),
 240 Multivariate Adaptive Regression Splines (MARS), Classification Tree Analysis (CTA), and
 241 eXtreme Gradient Boosting (XGBOOST) were selected to create an ensemble model for each
 242 butterfly species, using the ensemble platform “biomod2”⁴⁶ in R. We ensured that the
 243 underlying mechanism of our selection of algorithms was diverse and relatively balanced
 244 between the main categories of algorithms. We used 13 predictor variables for selection by
 245 individual models. All modeling was conducted at 10 km x 10 km resolution.

246

247 The Generalized Linear Model (GLM) is a regression-based algorithm widely used in
 248 SDMs⁴⁷. They are not as flexible when fitting complex response curve shapes, but this also
 249 means that GLMs are less vulnerable to overfitting⁴⁷. Maximum Entropy (MaxEnt) in our
 250 study was based on the “maxnet” R package⁴⁸, which uses penalized maximum likelihood for
 251 model fitting. MaxEnt is one of the computationally less expensive algorithms that perform
 252 well, making it a popular SDM algorithm⁴⁹. MaxEnt is more capable of fitting complicated,
 253 non-linear response curves, enabling users to model more complex relationships by using
 254 progressively complex statistics based on the number of samples available⁵⁰. The
 255 classification tree analysis (CTA) used by our SDM is based on the “rpart” R package⁵¹. The
 256 CTA algorithm recursively splits one group of data into two subgroups using one of the
 257 predictor variables given; therefore, the final model can be visualized as binary decision
 258 trees⁵¹. Finally, eXtreme Gradient Boosting (XGBoost) is one of the more computationally
 259 efficient gradient boosting algorithms implemented in R by the “xgboost” package⁵².
 260 Boosting algorithms feature an ensemble of weak models, each trained to minimize the errors
 261 of the previous models^{47,53}.

262

263 For the species distribution models, we used 13 predictor variables, which comprised 8
 264 Bioclim variables extracted from WorldClim⁵⁴, three soil variables extracted from SoilGrids⁵⁵
 265 through ISRIC (International Soil Reference and Information Centre)⁵⁶, and 2 vegetation
 266 variables derived from satellite data. The Bioclim variables employed included annual mean
 267 temperature (Bio 1), temperature seasonality (Bio 4), maximum temperature of warmest
 268 month (Bio 5), minimum temperature of coldest month (Bio 6), annual precipitation (Bio 12),
 269 precipitation of wettest month (Bio 13), precipitation of driest month (Bio 14), precipitation
 270 seasonality (Bio 15). The soil variables at a depth of 5-15 cm were used, including soil pH
 271 (phh2o), soil organic carbon content in the fine earth fraction (SOC), and total nitrogen
 272 (nitrogen). Nitrogen is generally recognized as one of the main limiting elements for plant
 273 growth⁵⁷, while soil organic carbon indicates soil quality⁵⁸. In addition, soil pH exerts
 274 considerable influence on soil biogeochemical processes, ultimately impacting plant
 275 growth⁵⁹. The selection of variables for our models was guided by expert knowledge to
 276 reflect/cover the key limitations and resources relevant to both butterflies and their host
 277 plants. Knowledge of the study region and biology/ecophysiology of the species being
 278 modeled allows the identification of the most ecologically relevant variables; therefore, it is
 279 the preferred approach for variable selection^{47,49,60,61}.

280

281 The vegetation variables used were the Normalized Difference Vegetation Index (NDVI) and
282 Canopy Height. NDVI was calculated from the USGS Landsat 5 (Level 2, Collection 2, Tier
283 1, 1985 – 1999) and USGS Landsat 7 (Level 2, Collection 2, Tier 1, 2000 – 2020) datasets,
284 with a customized script to filter satellite images by cloud cover (retaining images with 15%
285 or less cloud cover over land) and to obtain the mean NDVI value. Canopy Height data was
286 retrieved from the ETH Global Sentinel-2 10 m Canopy Height dataset⁶². These vegetation
287 cover variables were directly used to model the land cover/habitat available to butterflies.
288 Mean NDVI provided information on the general greenness of an area, while Canopy Height
289 data offered structural details on vegetation to better identify different types of habitats.
290 Together, these variables indicate resource availability and, to some extent, habitat structure.
291 To address potential issues associated with negative values in NDVI data, an alternative
292 variable, Corrected NDVI, which contains no negative values, was also examined. The
293 Corrected NDVI is derived from the equation $\text{Corrected NDVI} = \text{NDVI} + 1$. However, the
294 SDMs using Corrected NDVI produced identical results to those using standard NDVI data,
295 indicating that our models were unaffected by negative NDVI values.

296
297 The resolution of all environmental variables was set to 10 km x 10 km by averaging the
298 values from contributing grid cells. This resolution was chosen as a result of balancing the
299 spatial accuracy of available data and computational capabilities. Our data comprises 440,731
300 records with coordinate uncertainty data, while an additional 283,523 records that do not
301 have coordinate uncertainty data. Among the records with known coordinate uncertainty,
302 73,372 (18.31% of records with uncertainty data) had uncertainties ranging from 1-10 km,
303 and 39,302 (9.81% of records with uncertainty data) had uncertainties exceeding 10 km, thus
304 10 km seemed a reasonable compromise to reflect this. For the construction of SDMs, the
305 map of the study area and predicting variables were formatted to share the same extent,
306 resolution, and projection. We excluded entries with invalid species names (e.g., “NA” and
307 “not present”) or outside of our study area and those recorded before 1970. Data entries
308 published after 1970 but without date records were kept, assuming that their publication
309 infers their validity at the time of publication. Next, the map of tropical Asia and all
310 explanatory variable rasters were all projected to equal area projection EPSG:6933 and
311 cropped to the extent of 36° N to 10° S and 69° E to 161.6° E to fully cover the study region.
312 Our final cleaned database included 721,060 global records.

313
314 We used a function to further prepare the input files required by biomod2 and to generate
315 SDMs individually for each species. Occurrence data of a species was first extracted from our
316 butterfly occurrence dataset and used to produce a raster of resolution of 10 km x 10 km. A
317 total of n cells in the raster were assigned a value of “1” to represent at least one occurrence
318 record present in that cell, while cells with no record were assigned “n/a” instead of “0” since
319 no true absence data is available.

320
321 Only species with $n \geq 10$ were modeled. It has been shown that SDMs based on ten
322 occurrence points can reach 90% of the maximum possible accuracy⁶³, while recent studies
323 suggest a minimum requirement of 3 to 13 occurrence points in virtual simulations and 14 to
324 25 occurrence points in real-world conditions to infer accurate SDMs⁶⁴. Therefore, $n=10$ was
325 chosen as the lower limit of sample size for constructing SDMs to maximize the number of
326 species modeled while maintaining a reasonably high predictive accuracy⁶³. A total of 1,580
327 species met this qualification, whereas 1951 species had fewer records. For each species,
328 occurrence records were split into three sets: 10% of the data was first reserved for model
329 evaluation, and another 10% was then partitioned for model validation, leaving the remaining

330 80% of data for model calibration. The partitioning of model validation data was repeated 5
331 times to generate five different combinations of calibration and validation occurrence data.

332

333 Before SDM construction, pseudo-absence records were generated. Despite our efforts to fill
334 the spatial data gaps, the sampling effort of our dataset is still spatially biased toward highly
335 populated areas and roads due to the overwhelming number of records from GBIF and
336 iNaturalist in our dataset (more than 80%). As part of our effort to account for biases in our
337 data, we integrated the spatial bias of our dataset into the generation of pseudo-absence
338 records, assuming that all species were sampled in areas with at least one occurrence record
339 of any species. To capture such spatial bias, we created a raster layer of the spatial sampling
340 effort for all species across our study area (shown as sampling intensity in Fig. 1), which is
341 equivalent to the bias layer commonly used in the MaxEnt program. This was done by
342 pooling occurrence data of all species used in our models and summarising them in a raster,
343 then performing two-dimensional kernel density estimation (kde2d) using the R package
344 “MASS”⁶⁵ with the default settings. We excluded cells with occurrence records and sampled
345 the remaining study area for pseudo-absence records based on the bias layer, giving more
346 weight to well-sampled areas, as suggested by Phillips et al.⁶⁶ and Ferrier et al.⁶⁷. Following
347 the recommendation of Barbet-Massin et al.⁶⁸, for calibration, validation, and evaluation data,
348 we produced five sets of pseudo-absence data for each species, maintaining a 1:1 ratio
349 between the number of pseudo-absence points and occurrence points in each set.

350

351 Subsequently, we constructed SDMs for each species using five different partitions of
352 calibration and validation occurrence data, five selected algorithms, and five sets of pseudo-
353 absence data. This resulted in a total of 125 SDM models (5 x 5 x 5). Both presence and
354 pseudo-absence records were given equal weight during model construction to ensure a
355 consistent prevalence of 0.5 among all species. We applied a generalized setting for all
356 butterfly species for consistency across species, with adjustments made only to the learning
357 rate and the number of decision trees for the XGBoost algorithm to address overfitting. Other
358 model tuning options were retained at their default.

359

360 We generated binary outputs by maximizing True Skill Statistics (TSS), a widely used
361 threshold-dependent index of model fit. Ensemble modeling was selected over single best
362 models for its superior performance in rare species⁶⁹, and its robustness to uncertainties in
363 individual models by capturing the central tendency among models^{47,70,71}. We constructed an
364 ensemble model using all single models with TSS values greater than 0.7, ensuring that only
365 “substantial” models were included⁷². 1,520 species out of the 1,580 modeled species
366 obtained one or more single models meeting such criteria, allowing the further construction
367 of ensemble models. The ensemble model was generated using the mean algorithm⁷¹, where
368 all candidate models' probabilistic predictions were averaged without weighting. Finally, we
369 projected the ensemble model to the current environment using the same variables when
370 constructing the SDMs.

371

372 Ensemble models were evaluated using two metrics: TSS and Boyce index. TSS and Kappa
373 are two of the most popular SDM threshold-dependent evaluation metrics. TSS was chosen
374 over Kappa due to the inherent dependency of Kappa on species prevalence⁷³. Since we are
375 modeling thousands of species with differing degrees of rarity and prevalence, TSS is more
376 appropriate for model comparison between species. TSS varies from +1 to -1, in which +1
377 indicates perfect agreement with evaluation data, while a TSS value close to or less than 0
378 indicates model performance comparable to a random model⁷³.

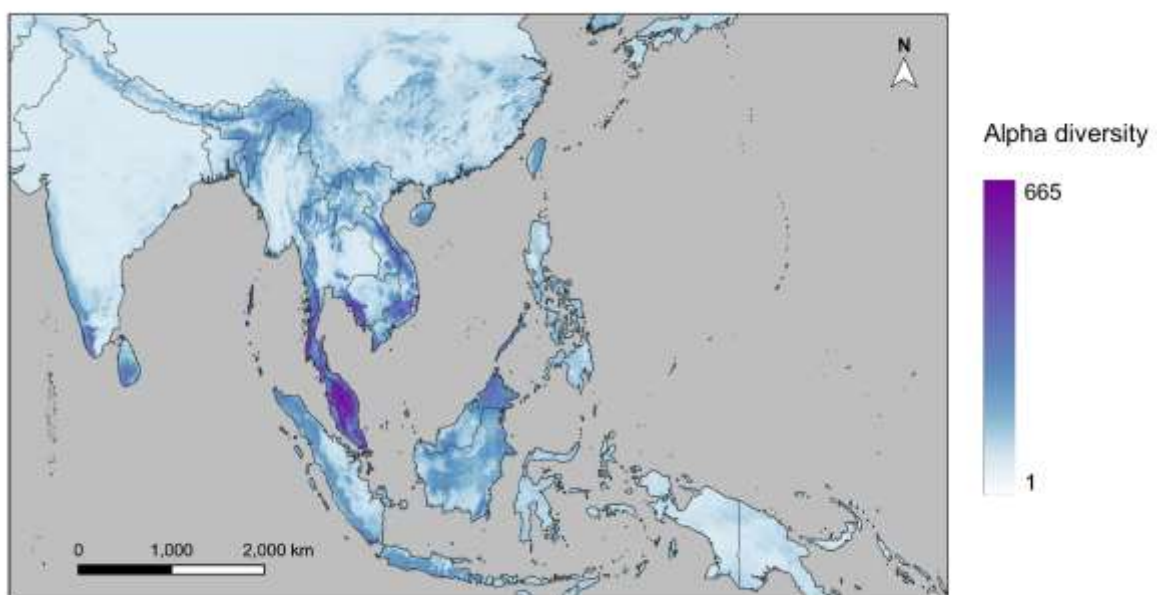
379

380 Following the suggestions of Hernandez et al.⁷⁴ and Breiner et al.⁶⁹ to use multiple evaluation
 381 measures when using presence-only data, we also calculated the Boyce index for all models
 382 built to supplement TSS. The Boyce index is capable of providing an accurate and reliable
 383 measure of model performance for models based on presence-only data⁷⁵, which is the key
 384 reason for its use in our study. Another reason for the use of the Boyce index is its lower
 385 sensitivity (correlation) to species prevalence relative to other metrics, including CVI,
 386 MaxKappa, and adjusted D2⁷⁵, while AUC and TSS also have a negative correlation with
 387 prevalence⁷³. AUC was also found to produce inflated estimates of model quality when the
 388 modeled species is rare⁷⁶. Boyce index ranges from +1 to -1, in which +1 indicates the model
 389 is of the highest quality and perfectly predicts evaluation data, while -1 indicates counter-
 390 prediction of evaluation data⁷⁵. Boyce index with a value close to 0 indicates the model
 391 performs no better than a random model⁷⁵.

392

393 To factor biogeography into predictions and correct for biogeographic overprediction
 394 generated by our SDMs (and account for differences between fundamental and realised
 395 niches), we restrained the sampling of pseudo-absence records and distribution maps
 396 produced by our models to regions that hosted more than 1% of species points (as such
 397 regions fall within species biogeographic ranges). By incorporating biogeography into model
 398 predictions, we aimed to reflect the impact of oceans as dispersal barriers in the SDM outputs
 399 to give a more realistic estimate of species' distribution and reduce false positive predictions.
 400 We first divided the landmasses of tropical Asia into 11 regions (Fig. 1) based on the
 401 ecoregions and biogeographic realms as revised by Dinerstein et al.⁴⁵, as well as Wallace's
 402 Line, Huxley's Line, and Weber's Line. For each species, we identified regions that included
 403 at least 1% of the species occurrence records, considering them to be active regions. We then
 404 cropped the SDM-predicted distribution maps to include only the active regions specific to
 405 each species. These cropped distribution maps were stacked together to generate an alpha
 406 diversity map, which illustrates the number of species present in each 10 km x 10 km cell
 407 across tropical Asia. The stacked SDM predictions highlighted a number of locations with
 408 relatively high diversity, exceeding 600 species in some locations (Fig. 2).

409



410

411 **Figure 2.** Projected distribution of butterfly diversity based on our species distribution
 412 models, using the mean algorithm for ensemble modeling.

413

414 *Point buffer methods*

415 For the 2,011 species (56% of all recorded species in our dataset) excluded from our species
 416 distribution modeling outputs either due to insufficient data or low quality of species
 417 distribution models, we plotted and buffered their occurrence records to infer alpha diversity.
 418 We first mapped their occurrence records and created 30 km-wide polygons (buffers) around
 419 these points to enhance clarity. Subsequently, the buffered occurrence points were converted
 420 into binary raster maps for each species and stacked to generate an additional alpha diversity
 421 map, representing species with limited occurrence records.

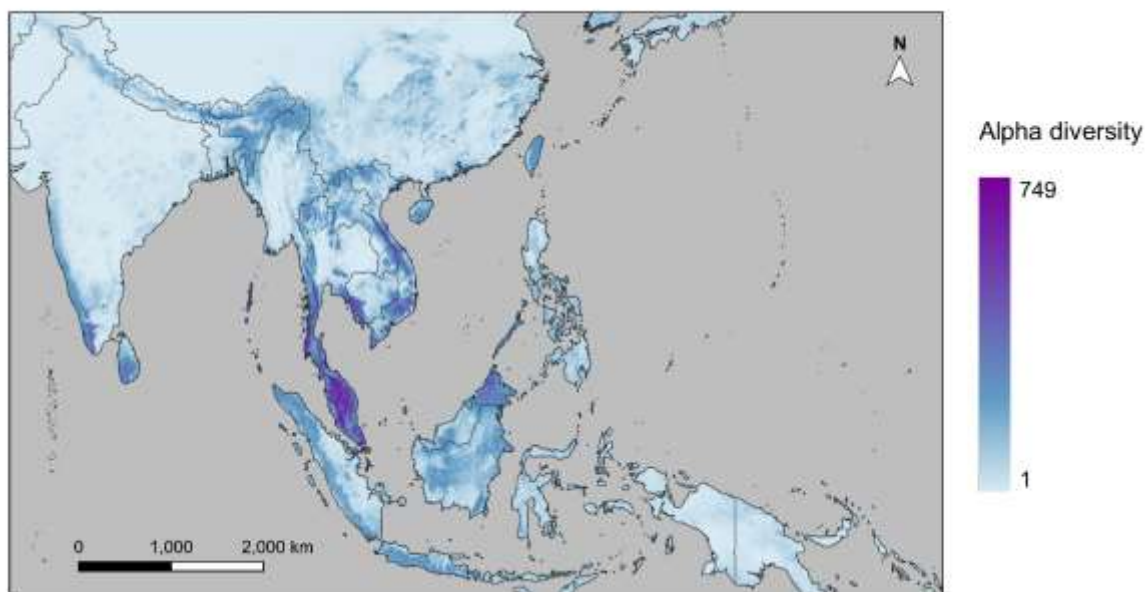
422

423 The diversity map derived from buffered occurrence points was then stacked with the species
 424 distribution model (SDM) projections to produce Fig. 3. This figure provides an overview of
 425 the alpha diversity of all species documented in our dataset. We identified two major
 426 butterfly diversity hotspots: peninsular Malaysia and the Sabah region of Borneo. We also
 427 found high levels of diversity predicted in Borneo, Sumatra, coastal Cambodia, southern
 428 Thailand, the Western Ghats in peninsular India, the Assam region of India, the Cardamom
 429 mountains in Cambodia, and Vietnam.

430

431

432



433

434 **Figure 3.** Estimated distribution of butterfly diversity based on our species distribution model
 435 projections and buffered occurrence points (for species not included in our SDM outputs).

436

437

438

439 *Software*

440 We calculated the SDMs in R⁷⁷, version 4.1.2. To construct and merge the SDMs into
 441 ensemble models, we utilized the "biomod2" package, version 4.2-4⁴⁶. The high-performance
 442 computing cluster HPC2021 at The University of Hong Kong, operating on CentOS 8, was
 443 employed to run the SDMs.

444

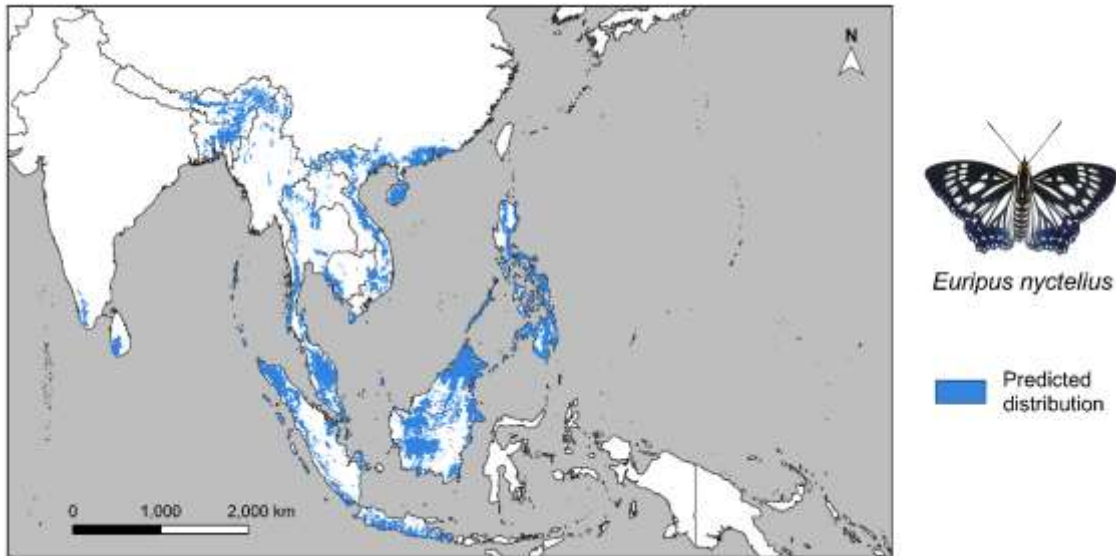
445

446 Data Records

447

448 All project files are publicly available in a Figshare repository
 449 (<https://doi.org/10.6084/m9.figshare.25037645>). Users may access the referenced occurrence
 450 dataset and metadata as .csv files, the SDM rasters (as file type, e.g., Fig. 4), and links to R
 451 scripts for SDM construction and distribution map generation were published on GitHub:
 452 <https://github.com/eugeneyau/Tropical-Asian-Butterfly-Distribution/tree/main>
 453 The GBIF-derived dataset is available at GBIF³⁷ (<https://doi.org/10.15468/dd.nvw5wr>).
 454 These outputs are licensed under a CC BY 4.0 license.

455



456

457 **Figure 4.** SDM-predicted distribution of *Euripus nyctelius* (Doubleday, 1845) (Nymphalidae:
 458 Apaturinae) based on our occurrence dataset.

459

460

461 Technical Validation

462

463 *SDM model evaluation/verification*

464 The mean TSS score of all ensemble models is 0.922, with a standard deviation of 0.155,
 465 while the Boyce index is 0.766, with a standard deviation of 0.305.

466

467 Both evaluation metrics indicate that the models constructed are of good quality. The mean
 468 TSS score of our ensemble models is higher than 0.8, falling into the category of “almost
 469 perfect” models according to the widely used division suggested by Landis & Koch⁷² (e.g.,
 470 Capinha et al.⁷⁸; Jones et al.⁷⁹). Since we only included models with TSS values of more than
 471 0.7 in our ensemble models, a high mean TSS score among the ensemble models is expected.
 472 The mean Boyce index of our models is higher than 0.7, which has been considered an
 473 indicator of good models in other studies (e.g., Rupprecht et al.⁸⁰). Boyce index value of 0.5
 474 is usually considered a cutoff for acceptable performance⁸¹.

475

476 *Collaborator Evaluation*

477 Our model outputs were also inspected by experts to evaluate their plausibility. Plausibility
 478 checks form an important part of model validation by making sure the modeling results

479 confine to the known range and possible range of the species modeled^{49,82}, serving as a
480 supplement to evaluation metrics, which only measure the goodness of fit of models.

481

482 Experts (coauthors/collaborators) agreed that our model outputs are generally reasonable and
483 informative. However, it is important to note that some of the sampling biases persisted in the
484 final model outputs despite our efforts to address data gaps by incorporating additional
485 datasets. We, therefore, encourage future data contributions to improve the coverage of our
486 dataset, especially in the areas with identified data gaps.

487

488 Although the majority of data gaps can be attributed to insufficient sampling effort, the
489 absence of data in the Philippines (and potentially other parts of tropical Asia) is primarily a
490 result of the dominance of Facebook over other platforms like iNaturalist for citizen science
491 data contribution. However, such data on Facebook contains limited information since EXIF
492 data (containing GPS coordinates) of photos are removed when uploaded. Filling the
493 Philippine data gap should be a priority, and mining Facebook data (e.g., Chowdhury et al.¹⁸)
494 and other sources would be a good place to start.

495

496 Our modeling results identified the Cardamom Mountains on the Cambodian-Thai border as a
497 butterfly diversity hotspot. During the Pleistocene when sea levels were up to 120 m lower
498 than present, and this area was on the eastern edge of a paleoriver watershed that included the
499 similarly diverse Malay peninsula and extended south to present-day Borneo^{83,84}. The high
500 diversity in this area is likely relictual⁸⁵. Endemism in this area likely contributes to high
501 butterfly diversity, which supports our models' prediction there.

502

503 Multiple experts pointed out the unexpected diversity differences between different parts of
504 Borneo. While our models identified Sabah as a hotspot for butterfly diversity, lower
505 diversity was predicted for other parts of Borneo, such as Sarawak and Kalimantan. This
506 contradicted our expectations, as all these areas possess mountainous regions and endemic
507 species, suggesting similar levels of butterfly diversity. The heart of Borneo, characterized by
508 lower disturbance compared to other parts of the island, was also predicted to host a relatively
509 lower diversity of butterflies by our models. Such a model prediction also contradicts our
510 expectation of higher butterfly diversity in less disturbed areas. This inconsistency between
511 expected and modeled butterfly diversity in Borneo is likely attributed to sampling bias,
512 evident through the alignment of modeled butterfly diversity with political boundaries and
513 sampling intensity (Fig. 1), and the lower modeled diversity in less accessible areas such as
514 the heart of Borneo (Fig. 2 and 3). The lack of data in less accessible areas has been discussed
515 by Hughes et al.¹⁵ and Boakes et al.⁸⁶, while this trend is even more obvious in citizen science
516 data¹⁵, which constitutes a considerable proportion of our dataset.

517

518 While some of the spatial variations in the sampling effort of our dataset are reflected in the
519 spatial bias of our modeling results, there are several notable discrepancies between the
520 distribution of data and modeled diversity. Fig. 1 illustrates that Japan, Taiwan, and northern
521 Thailand have a relatively high intensity of sampling effort compared to their predicted
522 butterfly diversity in Fig. 2. Conversely, a reversed pattern is evident in Southern Borneo and
523 Southern Sumatra, where our data shows low sampling effort but our models predict high
524 butterfly diversity. These patterns demonstrate the robustness of the models to some of the
525 spatial sampling biases present in our data.

526

527 To determine the variable importance in our SDMs, we calculated, for each variable, the
528 mean variable importance throughout the ensemble models of all species. Temperature

529 seasonality (Bio 4) emerged as the most important variable (scoring 0.280 out of 1), followed
 530 by the minimum temperature of the coldest month (Bio 6, scoring 0.163 out of 1) and annual
 531 mean temperature (Bio 1, scoring 0.140 out of 1). Soil pH (pH2o, scoring 0.107 out of 1),
 532 precipitation of driest month (Bio 14, scoring 0.0973 out of 1), and Canopy Height (scoring
 533 0.0907 out of 1) also exhibited high importance in the models. The ranking of variable
 534 importance in the SDMs conforms to the hierarchical framework of Pearson & Dawson⁸⁷, in
 535 which climatic variables exert greater control over species distribution at continental scales,
 536 while land cover and soil variables gain influence at more localized scales. In addition, the
 537 high importance of temperature variables, particularly temperature seasonality (Bio 4), is
 538 consistent with the results of Carvalho et al.⁸⁸, which highlighted the strong impact of
 539 temperature, especially temperature seasonality, on butterfly distribution and diversity.

540

541

542 Usage Notes

543

544 The predictor variables considered in our SDMs, which include the 8 Bioclim variables and
 545 the 3 SoilGrids variables, are products of interpolation between available point data^{54,55}. As
 546 with most data collected without stratified sampling, these point data are likely to be spatially
 547 biased. Users should note that our SDMs inherit these biases, as well as uncertainties in the
 548 interpolation result.

549

550 By generating more pseudo-absences for SDMs in well-sampled areas with the use of the bias
 551 mask, we are essentially augmenting the weighting of extensively surveyed regions in our
 552 models, while unsampled habitats may be presumed as suitable. Consequently, the
 553 transferability of our models to unsampled areas is limited, especially when extrapolating in
 554 novel environments not covered by training data⁶⁶ or in areas where biogeographic barriers
 555 prevent dispersal. This is also one of the reasons for restraining our model predictions to the
 556 regions where a species is known to occur so that the results are not overly optimistic. Such
 557 an approach to pseudo-absence generation also assumes that the data collection method is
 558 consistent throughout the entire dataset⁶⁶, while our dataset is compiled from various sources.
 559 However, since a majority of our data is derived from a single source (GBIF), we can
 560 consider the data collection method consistent in terms of the observation method. To use our
 561 data and models for the prediction of future butterfly distribution under climate change, we
 562 suggest using the “random” method from the biomod2 package to generate pseudo-absence
 563 records.

564

565 Regarding uncertainty in model results, we have limited confidence in the model predictions
 566 in the Philippines and New Guinea. The scarcity of occurrence data in these two regions (Fig.
 567 1) prevents strong inferences. Additionally, the presence of biogeographic barriers such as
 568 Wallace's Line and Huxley's Line restricts the use of occurrence data from other regions to
 569 infer butterfly distribution in these specific areas.

570

571

572 Code Availability

573

574 All code used to conduct synonym harmonization, preprocess environmental variables for
 575 SDMs, execute SDMs, process SDM outputs, and conduct point buffer analysis can be
 576 accessed in our GitHub project repository:

577 <https://github.com/eugeneyau/Tropical-Asian-Butterfly-Distribution>

578 under the Code directory. All data, including our butterfly occurrence dataset, SDM-predicted
579 distribution maps, tropical Asian biogeographic regions shapefile, and buffered occurrence
580 points for species excluded from our species distribution modeling are available from our
581 FigShare repository (<https://doi.org/10.6084/m9.figshare.25037645>). Some data for Sulawesi,
582 provided by TK, have not been included in the data release but are available upon request to
583 TK.

584
585

586 **Acknowledgments**

587

588 Funding for this research was provided by a National Science Foundation China Excellent
589 Young Scientist award to TCB. Additionally, preliminary data and analyses were supported
590 by a Global Environmental Facility grant. The computations were performed using research
591 computing facilities offered by Information Technology Services at The University of Hong
592 Kong. Landsat-5 and Landsat-7 image courtesy of the U.S. Geological Survey. Special thanks
593 to Kirsten Boehm, Ryan Leung, Xueying Wang, and Tracy Zhang for assistance with data
594 extraction. We are very grateful to the Natural History Museum UK (NHMUK) and
595 Zoologische Staatssammlung München (ZSM) for facilitating TK's access to the collection
596 and would like to thank the curators for their kind support. We wish to thank B. Huertas, C.
597 Beale, and R. Vane-Wright for their constant support of TK's work. LVV was supported by
598 the Vietnam Ministry of Science and Technology (ĐTĐL.CN-113/21).

599
600

601 **Author Contributions**

602

603 TCB, TPNT, SX, RTC, and PR conceptualized the study, with funding acquired by TCB,
604 TPNT, and SX. The project methodology was developed by EEJ, EYHY, TPNT, SX, AKCL,
605 RTC, PR, and ACH. TCB and ACH supervised the study, while TCB managed the project
606 administration. The occurrence dataset was compiled by EEJ with assistance from CWCH,
607 and data curation was carried out by EEJ, EYHY, and DJL. Data were contributed by DJL,
608 SC, JKH, YB, TK, LVV, and AN. EYHY conducted the species distribution modeling and
609 subsequent analyses. DJL validated scientific names in the dataset and, along with CWCH,
610 SC, JKH, ALM, JAT, ICC, GC, GCW, SB, AJ, TK, KK, and SAS, provided insights on the
611 plausibility of the species distribution model outputs. The initial draft of the manuscript was
612 written by EEJ, EYHY, TCB, and ACH, with EYHY responsible for data visualization. All
613 authors contributed input and suggestions to the draft and approved the final manuscript.

614
615

616 **Additional information**

617
618

619 **Competing interests**

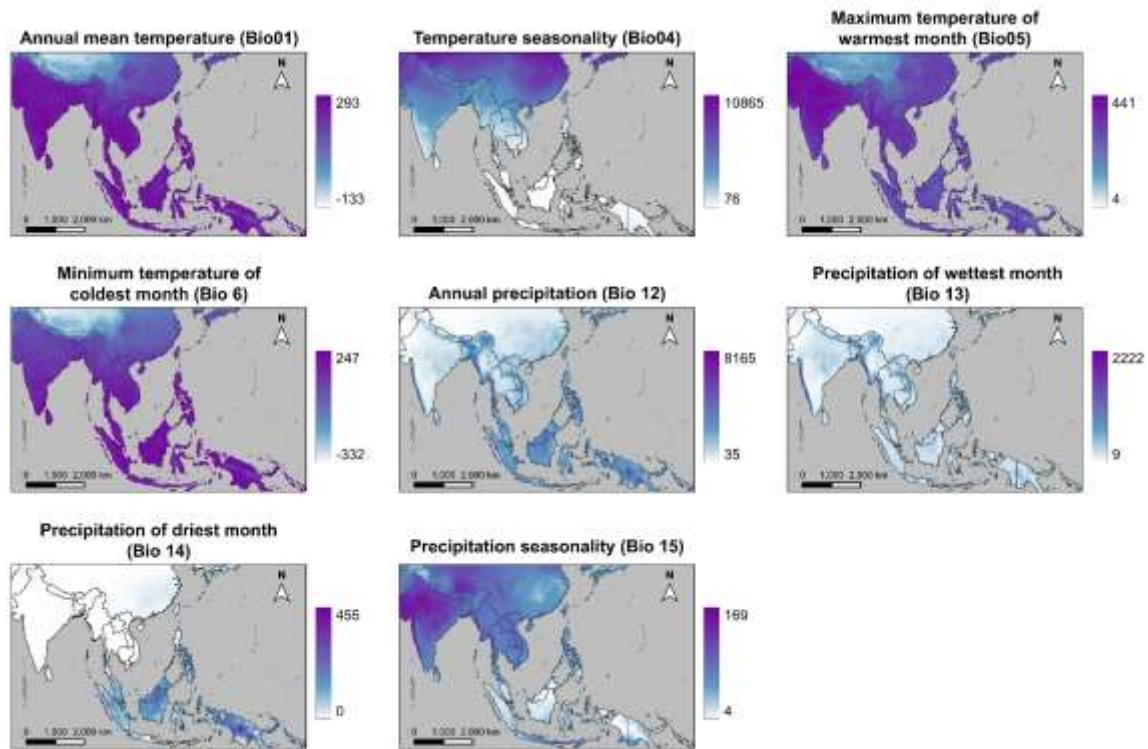
620

621 The authors declare no competing interests.

622
623

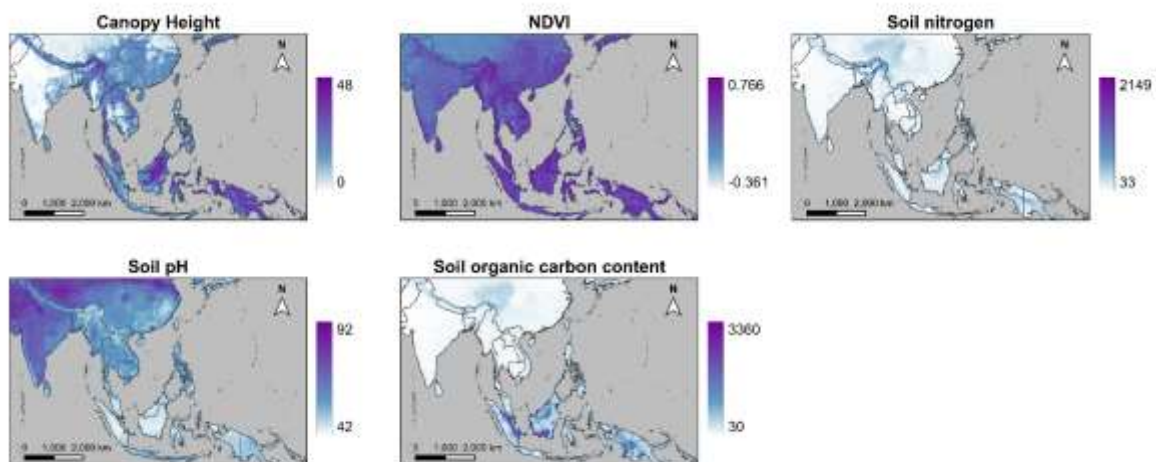
624 **Supplementary Information**

625



626
627
628

Supplementary Figure 1. Climatic predictor variables included in our SDMs.



629
630
631
632

Supplementary Figure 2. Non-climatic predictor variables included in our SDMs.

References

633
634
635
636
637
638
639
640
641

1. de Bruyn, M. et al. Borneo and Indochina are major evolutionary hotspots for Southeast Asian biodiversity. *Syst. Biol.* **63**, 879-901 (2014).
2. Verde Arregoitia, L. D. Biases, gaps, and opportunities in mammalian extinction risk research. *Mammal Rev.* **46**, 17-29 (2016).
3. Cardoso, P., Erwin, T. L., Borges, P. A. & New, T. R. The seven impediments in invertebrate conservation and how to overcome them. *Biol. Conserv.* **144**, 2647-2655 (2011).

- 642 4. Xing, Shuang. et al. Conservation of data deficient species under multiple threats:
643 Lessons from an iconic tropical butterfly (*Teinopalpus aureus*). *Biol. Conserv.* **234**, 154-
644 164 (2019).
- 645 5. Pinkert, S., Barve, V., Guralnick, R. & Jetz, W. Global geographical and latitudinal
646 variation in butterfly species richness captured through a comprehensive country-level
647 occurrence database. *Glob. Ecol. Biogeogr.* **31**, 830-839 (2022).
- 648 6. Kawahara, A. Y. et al. A global phylogeny of butterflies reveals their evolutionary
649 history, ancestral hosts and biogeographic origins. *Nat. Ecol. Evol.* **7**, 903–913 (2023).
- 650 7. Pinkert, S., Sica, Y. V., Winner, K. & Jetz, W. The potential of ecoregional range maps
651 for boosting taxonomic coverage in ecology and conservation. *Ecography* **2023**, (2023).
- 652 8. Whittaker, R. J. et al. Conservation biogeography: assessment and prospect. *Divers.*
653 *Distrib.* **11**, 3-23 (2005).
- 654 9. Earl, C. et al. Spatial phylogenetics of butterflies in relation to environmental drivers and
655 angiosperm diversity across North America. *IScience* **24**, (2021).
- 656 10. Wei, C.-H., Lohman, D. J., Peggie, D. & Yen, S.-H. An illustrated checklist of the genus
657 *Elymnias* Hübner, 1818 (Nymphalidae, Satyrinae). *Zookeys* **676**, 47-152. (2017).
- 658 11. Condamine, F. L. et al. Fine-scale biogeographical and temporal diversification
659 processes of peacock swallowtails (*Papilio* subgenus *Achillides*) in the Indo-Australian
660 Archipelago. *Cladistics* **29**, 88–111 (2012).
- 661 12. Toussaint, E. F. et al. Comparative molecular species delimitation in the charismatic
662 Nawab butterflies (Nymphalidae, Charaxinae, *Polyura*). *Mol. Phylogenet. Evol.* **91**, 194-
663 209 (2015).
- 664 13. Scriven, S. A. et al. Assessing the effectiveness of protected areas for conserving range-
665 restricted rain forest butterflies in Sabah, Borneo. *Biotropica* **52**, 380-391 (2020).
- 666 14. Beck, J., Böller, M., Erhardt, A. & Schwanghart, W. Spatial bias in the GBIF database
667 and its effect on modeling species' geographic distributions. *Ecol. Inform.* **19**, 10-15
668 (2014).
- 669 15. Hughes, A. C. et al. Sampling biases shape our view of the natural world. *Ecography* **44**,
670 1259-1269 (2021).
- 671 16. Meyer, C., Kreft, H., Guralnick, R. & Jetz, W. Global priorities for an effective
672 information basis of biodiversity distributions. *Nat. Commun.* **6**, 1-8 (2015).
- 673 17. Peterson, A. T., Navarro-Sigüenza, A. G. & Gordillo, A. Assumption-versus data-based
674 approaches to summarizing species' ranges. *Conserv. Biol.* **32**, 568-575 (2016).
- 675 18. Chowdhury, S. et al. Using social media records to inform conservation planning.
676 *Conserv. Biol.* **38**, (2024).
- 677 19. Jetz, W., Sekercioglu, C. H. & Watson, J. E. Ecological correlates and conservation
678 implications of overestimating species geographic ranges. *Conserv. Biol.* **22**, 110-119
679 (2008).
- 680 20. Sofaer, H. R. et al. Development and delivery of species distribution models to inform
681 decision-making. *Biosci.* **69**, 544-557 (2019).
- 682 21. Anderson, R. P. et al. *Ecological niches and geographic distributions* (Princeton
683 University Press, 2011).
- 684 22. Troudet, J., Grandcolas, P., Blin, A., Vignes-Lebbe, R. & Legendre, F. Taxonomic bias
685 in biodiversity data and societal preferences. *Sci. Rep.* **7**, (2017).
- 686 23. Kramer-Schadt, S. et al. The importance of correcting for sampling bias in MaxEnt
687 species distribution models. *Divers. Distrib.* **19**, 1366-1379 (2013).
- 688 24. Guillera-Aroita, G. et al. Is my species distribution model fit for purpose? Matching
689 data and models to applications. *Glob. Ecol. Biogeogr.* **24**, 276-292 (2015).

- 690 25. Smeraldo, S. et al. Species distribution models as a tool to predict range expansion after
691 reintroduction: A case study on Eurasian beavers (*Castor fiber*). *J. Nat. Conserv.* **37**, 12-
692 20 (2017).
- 693 26. Guisan, A. et al. Predicting species distributions for conservation decisions. *Ecol. Lett.*
694 **16**, 1424-1435 (2013).
- 695 27. Araújo, M. B. et al. Standards for distribution models in biodiversity assessments. *Sci.*
696 *Adv.* **5**, (2019).
- 697 28. Maes, D. et al. The potential of species distribution modelling for reintroduction
698 projects: the case study of the Chequered Skipper in England. *J. Insect Conserv.* **23**,
699 419-431 (2019).
- 700 29. Pacifici, M. et al. Assessing species vulnerability to climate change. *Nat. Clim. Change*
701 **5**, 215-224 (2015).
- 702 30. Attorre, F. et al. How to include the impact of climate change in the extinction risk
703 assessment of policy plant species? *J. Nat. Conserv.* **44**, 43-49 (2018).
- 704 31. Lohman, D. J. et al. Biogeography of the Indo-Australian archipelago. *Annu. Rev. Ecol.*
705 *Evol. Syst.* **42**, 205-226 (2011).
- 706 32. Hughes, A. C. Understanding the drivers of Southeast Asian biodiversity loss.
707 *Ecosphere* **8**, (2017).
- 708 33. Corlett, R. T. *The Ecology of Tropical East Asia* 3rd edn (Oxford University Press,
709 2019).
- 710 34. Sodhi, N. S., Koh, L. P., Brook, B. W. & Ng, P. K. L. Southeast Asian biodiversity: an
711 impending disaster. *Trends Ecol. Evol.* **19**, 654-660 (2004).
- 712 35. Wilcove, D. S., Giam, X., Edwards, D. P., Fisher, B., & Koh, L. P. Navjot's nightmare
713 revisited: logging, agriculture, and biodiversity in Southeast Asia. *Trends Ecol. Evol.* **28**,
714 531-540 (2013).
- 715 36. Orr, M. C. et al. Global patterns and drivers of bee distribution. *Curr Biol.* **31**, 451-458
716 (2021).
- 717 37. *Global Biodiversity Information Facility* <https://doi.org/10.15468/dd.nvw5wr> (2024).
- 718 38. Ball-Damerow, J. E. et al. Research applications of primary biodiversity databases in the
719 digital age. *PloS one* **14**, (2019).
- 720 39. Gaiji, S. et al. Content assessment of the primary biodiversity data published through
721 GBIF network: status, challenges and potentials. *Biodiversity Informatics* **8**, (2013).
- 722 40. Costello, M. J., Michener, W. K., Gahegan, M., Zhang, Z. Q. & Bourne, P. E.
723 Biodiversity data should be published, cited, and peer reviewed. *Trends Ecol. Evol.* **28**,
724 454-461 (2013).
- 725 41. Goodwin, Z. A., Harris, D. J., Filer, D., Wood, J. R. & Scotland, R. W. Widespread
726 mistaken identity in tropical plant collections. *Curr. Biol.* **25**, R1066-R1067 (2015).
- 727 42. Chowdhury, S. et al. Butterflies are weakly protected in a mega-populated country,
728 Bangladesh. *Glob. Ecol. Conserv.* **26**, (2021).
- 729 43. Shields, O. World numbers of butterflies. *J. Lepid. Soc.* **43**, 178-183 (1989).
- 730 44. Robbins, R. K. & Opler, P. A. in *Biodiversity II: Understanding and Protecting Our*
731 *Biological Resources* (eds. Reaka-Kudla, M. L., Wilson, D. E. & Wilson, E. O.) Ch. 6
732 (Joseph Henry Press, 1997).
- 733 45. Dinerstein, E. et al. An ecoregion-based approach to protecting half the terrestrial realm.
734 *Biosci.* **67**, 534-545 (2017).
- 735 46. Thuiller, W. et al. biomod2: Ensemble Platform for Species Distribution Modeling. R
736 package version 4.2-4. <https://CRAN.R-project.org/package=biomod2> (2023).
- 737 47. Guisan, A., Thuiller, W. & Zimmermann, N. E. *Habitat Suitability and Distribution*
738 *Models: With Applications in R* (Cambridge University Press, 2017).

- 739 48. Phillips, S. J. maxnet: Fitting ‘Maxent’ Species Distribution Models with ‘glmnet’. R
740 package version 0.1.4. <https://CRAN.R-project.org/package=maxnet> (2021).
- 741 49. Porfirio, L. L. et al. Improving the use of species distribution models in conservation
742 planning and management under climate change. *PloS One* **9**, (2014).
- 743 50. Elith, J. & Graham, C. H. Do they? How do they? WHY do they differ? On finding
744 reasons for differing performances of species distribution models. *Ecography* **32**, 66-77
745 (2009).
- 746 51. Therneau, T., Atkinson, B. & Ripley, B. Rpart: Recursive Partitioning. R Package
747 version 4.1-3. <http://CRAN.R-project.org/package=rpart> (2013).
- 748 52. Chen, T. et al. xgboost: Extreme Gradient Boosting. R package version 1.7.5.1.
749 <https://CRAN.R-project.org/package=xgboost> (2023).
- 750 53. Natekin, A. & Knoll, A. Gradient boosting machines, a tutorial. *Front. Neurobot.* **7**,
751 (2013).
- 752 54. Hijmans, R. J., Cameron, S. E., Parra, J. L., Jones, P. G. & Jarvis, A. Very high
753 resolution interpolated climate surfaces for global land areas. *Int. J. Climatol.* **25**, 1965-
754 1978 (2005).
- 755 55. Poggio, L. et al. SoilGrids 2.0: producing soil information for the globe with quantified
756 spatial uncertainty. *SOIL* **7**, 217-240 (2021).
- 757 56. *International Soil Reference and Information Centre* <https://www.isric.org/> (2024).
- 758 57. Ågren, G. I., Wetterstedt, J. M. & Billberger, M. F. Nutrient limitation on terrestrial
759 plant growth—modelling the interaction between nitrogen and phosphorus. *New Phytol.*
760 **194**, 953-960 (2012).
- 761 58. Lal, R. Soil health and carbon management. *Food Energy Secur.* **5**, 212-222 (2016).
- 762 59. Neina, D. The role of soil pH in plant nutrition and soil remediation. *Appl. Environ. Soil*
763 *Sci.* **2019**, 1-9 (2019).
- 764 60. Barbet-Massin, M. & Jetz, W. A 40-year, continent-wide, multispecies assessment of
765 relevant climate predictors for species distribution modelling. *Divers. Distrib.* **20**, 1285-
766 1295 (2014).
- 767 61. Zeng, Y., Low, B. W. & Yeo, D. C. Novel methods to select environmental variables in
768 MaxEnt: A case study using invasive crayfish. *Ecol. Model.* **341**, 5-13 (2016).
- 769 62. Lang, N., Jetz, W., Schindler, K. & Wegner, J. D. A high-resolution canopy height
770 model of the Earth. *Nat. Ecol. Evol.* **7**, 1778-1789 (2023).
- 771 63. Stockwell, D. R. & Peterson, A. T. Effects of sample size on accuracy of species
772 distribution models. *Ecol. Model.* **148**, 1-13 (2002).
- 773 64. van Proosdij, A. S., Sosef, M. S., Wieringa, J. J. & Raes, N. Minimum required number
774 of specimen records to develop accurate species distribution models. *Ecography* **39**,
775 542-552 (2016).
- 776 65. Venables, W. N. & Ripley, B. D. *Modern Applied Statistics with S-PLUS* 4th edn
777 (Springer, 2002).
- 778 66. Phillips, S. J. et al. Sample selection bias and presence-only distribution models:
779 implications for background and pseudo-absence data. *Ecol. Appl.* **19**, 181-197 (2009).
- 780 67. Ferrier, S., Watson, G., Pearce, J. & Drielsma, M. Extended statistical approaches to
781 modelling spatial pattern in biodiversity in northeast New South Wales. I. Species-level
782 modelling. *Biodiversity Conserv.* **11**, 2275-2307 (2002).
- 783 68. Barbet-Massin, M., Jiguet, F., Albert, C. H. & Thuiller, W. Selecting pseudo-absences
784 for species distribution models: How, where and how many?. *Methods Ecol. Evol.* **3**,
785 327-338 (2012).
- 786 69. Breiner, F. T., Guisan, A., Bergamini, A. & Nobis, M. P. Overcoming limitations of
787 modelling rare species by using ensembles of small models. *Methods Ecol. Evol.* **6**,
788 1210-1218 (2015).

- 789 70. Araújo, M. B. & New, M. Ensemble forecasting of species distributions. *Trends Ecol.*
790 *Evol.* **22**, 42-47 (2007).
- 791 71. Marmion, M., Parviainen, M., Luoto, M., Heikkinen, R. K. & Thuiller, W. Evaluation of
792 consensus methods in predictive species distribution modelling. *Divers. Distrib.* **15**, 59-
793 69 (2009).
- 794 72. Landis, J. R. & Koch, G. G. The measurement of observer agreement for categorical
795 data. *Biometrics* **33**, 159-174 (1977).
- 796 73. Allouche, O., Tsoar, A. & Kadmon, R. Assessing the accuracy of species distribution
797 models: prevalence, kappa and the true skill statistic (TSS). *J. Appl. Ecol.* **43**, 1223-1232
798 (2006).
- 799 74. Hernandez, P. A., Graham, C. H., Master, L. L. & Albert, D. L. The effect of sample
800 size and species characteristics on performance of different species distribution
801 modeling methods. *Ecography* **29**, 773-785 (2006).
- 802 75. Hirzel, A. H., Le Lay, G., Helfer, V., Randin, C. & Guisan, A. Evaluating the ability of
803 habitat suitability models to predict species presences. *Ecol. Model.* **199**, 142-152
804 (2006).
- 805 76. Lobo, J. M., Jiménez-Valverde, A. & Real, R. AUC: a misleading measure of the
806 performance of predictive distribution models. *Glob. Ecol. Biogeogr.* **17**, 145-151
807 (2008).
- 808 77. R Core Team. R: A language and environment for statistical computing, version 4.1.2. R
809 Foundation for Statistical Computing <https://www.R-project.org/> (2021).
- 810 78. Capinha, C., Rocha, J. & Sousa, C. A. Macroclimate determines the global range limit of
811 *Aedes aegypti*. *Ecohealth* **11**, 420-428 (2014).
- 812 79. Jones, C. C., Acker, S. A. & Halpern, C. B. Combining local-and large-scale models to
813 predict the distributions of invasive plant species. *Ecol. Appl.* **20**, 311-326 (2010).
- 814 80. Rupprecht, F., Oldeland, J. & Finckh, M. Modelling potential distribution of the
815 threatened tree species *Juniperus oxycedrus*: how to evaluate the predictions of different
816 modelling approaches? *J. Veg. Sci.* **22**, 647-659 (2011).
- 817 81. Pomoim, N., Hughes, A. C., Trisurat, Y. & Corlett, R. T. Vulnerability to climate
818 change of species in protected areas in Thailand. *Sci. Rep.* **12**, (2022).
- 819 82. Zurell, D. A standard protocol for reporting species distribution models. *Ecography* **43**,
820 1261-1277 (2020).
- 821 83. Sholihah, A. et al. Impact of Pleistocene eustatic fluctuations on evolutionary dynamics
822 in Southeast Asian biodiversity hotspots. *Syst. Biol.* **70**, 940-960 (2021).
- 823 84. Voris, H. K. Maps of Pleistocene sea levels in Southeast Asia: Shorelines, river systems
824 and time durations. *J. Biogeogr.* **27**, 1153-1167 (2000).
- 825 85. Monastyrskii, A. L. & Vane-Wright, R. I. Identity of *Euploea orontobates* Fruhstorfer,
826 1910 (Lepidoptera: Nymphalidae), a milkweed butterfly from Thailand and Vietnam.
827 *Zootaxa* **1991**, 43-50 (2009).
- 828 86. Boakes, E. H. et al. Distorted views of biodiversity: spatial and temporal bias in species
829 occurrence data. *PLoS Biol.* **8**, (2010).
- 830 87. Pearson, R. G. & Dawson, T. P. Predicting the impacts of climate change on the
831 distribution of species: are bioclimate envelope models useful? *Glob. Ecol. Biogeogr.*
832 **12**, 361-371 (2003).
- 833 88. Carvalho, A. P. S. et al. Comprehensive phylogeny of Pieridae butterflies reveals strong
834 correlation between diversification and temperature. *iScience* **27**, (2024).
- 835