

1 **WALLACE: A FLEXIBLE PLATFORM FOR REPRODUCIBLE MODELING OF SPECIES NICHES**  
2 **AND DISTRIBUTIONS BUILT FOR COMMUNITY EXPANSION**

3

4 **JAMIE M. KASS, BRUNO VILELA, MATTHEW E. AIELLO-LAMMENS, ROBERT**  
5 **MUSCARELLA, CORY MEROW AND ROBERT P. ANDERSON**

6

7 **ABTRACT**

- 8 1) Scientific research increasingly calls for open-source software that is flexible,  
9 interactive, and expandable, while providing methodological guidance and  
10 reproducibility. Currently, many analyses in ecology are implemented with “black box”  
11 graphical user interfaces that lack flexibility or command-line interfaces that are  
12 infrequently used by non-specialists.
- 13 2) To help remedy this situation in the context of species distribution modeling, we  
14 created *Wallace*, an open and modular application with a richly documented graphical  
15 user interface to underlying R scripts that is flexible and highly interactive.
- 16 3) *Wallace* guides users from acquiring and processing data to building models and  
17 examining predictions. Additionally, it is designed to grow via community contributions  
18 of new modules to expand functionality. All results are downloadable, along with code to  
19 reproduce the analysis.
- 20 4) *Wallace* provides an example of an innovative platform to increase access to cutting-  
21 edge methods and encourage plurality in science and collaboration in software  
22 development.

## 23 INTRODUCTION

24 Ecological and evolutionary studies have shifted over the past 20 years towards  
25 increasingly complex analyses (Bolker 2008). This has been enabled, in part, by a rise in  
26 computing power and the increasing openness of data and software (Gimenez *et al.* 2014).  
27 As a result, most current methods are accessed as either: 1) programming-language  
28 scripts run in command-line interfaces (CLIs; e.g. R and Python), or 2) software with  
29 graphical user interfaces (GUIs). On one hand, programming scripts provide flexibility,  
30 but custom code is often poorly documented and tailored to specific analyses (Mislan *et*  
31 *al.* 2016). GUIs, on the other hand, are easy to navigate and extend accessibility of  
32 analyses to more users, but are less flexible than custom code and often necessitate using  
33 multiple software packages to complete a study. This exacerbates a problem with GUIs:  
34 lack of reproducibility (Hampton *et al.* 2015). Additionally, GUI implementations of  
35 methods often lag behind the cutting-edge analyses enabled by the frequent release of  
36 scripts with new publications. Hence, tools that combine the positive aspects of CLI and  
37 GUI methods can help advance ecological research.

38

39 We developed *Wallace* to address these issues specifically for user communities in  
40 ecology and the environmental sciences. *Wallace* is an open-source GUI application that  
41 offers user-friendly access to R-scripted modern workflows. It is available as the R  
42 package “wallace” on CRAN, with a development version on Github (see Data  
43 Accessibility for links). *Wallace* currently focuses on a workflow for modeling species  
44 niches and geographic distributions (Fig. 1; Guisan & Thuiller 2005; Peterson *et al.*  
45 2011), but we anticipate that future versions will expand the analyses offered for

46 biogeographical and macroecological modeling. *Wallace* is written for R (R Core Team  
47 2017) using *shiny* (a package for developing interactive applications; Chang *et al.*  
48 2017), and can thus leverage the rapidly expanding suite of R packages authored by the  
49 scientific community. Six main qualities of *Wallace* distinguish it as a model for  
50 providing access and guidance for advanced methodologies (Table 1).

51

52 \*Figure 1\*

53

54 Below, we present several important issues in niche/distribution modeling and explain  
55 how we address them with *Wallace*, first conceived as a response to the Global  
56 Biodiversity Information Facility’s 2015 Ebbe Nielsen Challenge  
57 (<https://devpost.com/software/wallace-round-2>). We then provide a walkthrough of the  
58 application and conclude by discussing the general utility of *Wallace*’s framework for  
59 disseminating scientific methods and encouraging community-wide innovation.

60

61 \*Table 1\*

62

### 63 **CURRENT ISSUES IN NICHE/DISTRIBUTION MODELING**

64 *Wallace* currently implements analyses for species niche/distribution modeling (hereafter  
65 “distribution modeling”). These correlative models estimate the response of a species to  
66 the environment and with clear assumptions can be used to infer (or hypothesize)  
67 geographic ranges, environmental suitability across a landscape, or niche requirements  
68 (Franklin 2010a; Peterson *et al.* 2011). Distribution modeling is used in many disciplines,

69 such as phylogeography (Alvarado-Serrano & Knowles 2014), community ecology  
70 (Guisan & Rahbek 2011), evolutionary biology (McCormack *et al.* 2010), and  
71 conservation (Franklin 2010b). At a minimum, it requires georeferenced occurrence  
72 records of the study species (e.g., from field surveys, museum collections, citizen  
73 science) and environmental predictors (e.g., climate, land cover, topography). Occurrence  
74 data generally represent the primary ecological information available for the vast  
75 majority of species. “Presence-only” distribution models use environmental values at  
76 occurrences, typically contrasting them with those available in the study region  
77 (“background” or pseudoabsence samples; Elith *et al.* 2011). Since reliable absence data  
78 are unavailable for most species, research focusing on presence-only models has grown  
79 tremendously over the past two decades. The current implementation of *Wallace*  
80 concentrates on these models, highlighting two algorithms with differing complexity:  
81 BIOCLIM (Booth *et al.* 2014) and Maxent (Phillips *et al.* 2006). Many approaches exist  
82 for making such models, and comparing them conveys to users that the utility of a model  
83 does not necessarily improve with complexity (Jiménez-Valverde *et al.* 2008).

84       Confusion abounds regarding how best to choose and implement presence-only  
85 distribution modeling methods and interpret their outputs (Joppa *et al.* 2013). There have  
86 been numerous calls to address a range of complicating issues (Elith *et al.* 2011; Merow  
87 *et al.* 2013), among them sampling bias (Bean *et al.* 2012), selection of study extent  
88 (VanDerWal *et al.* 2009), model evaluation (Radosavljevic & Anderson 2014), model  
89 selection (Warren & Seifert 2011), and considering key assumptions (Yackulic *et al.*  
90 2013). *Wallace* provides extensive guidance text and enables user experimentation with a

91 variety of modules (see *Walkthrough*), directly addressing some of these issues and  
92 encouraging the use of a diversity of methods (Fig. 2).

93

94 \*Figure 2\*

95

## 96 **DIFFICULTIES FOR PROGRAMMERS AND NON-PROGRAMMERS ALIKE**

97 *Wallace* combines the strengths of GUI and CLI approaches to enable research in  
98 distribution modeling for a broad audience. A number of GUI-based applications have  
99 been widely used for distribution modeling analyses (e.g., maxent.jar (Phillips *et al.*  
100 2006); DesktopGARP (Scachetti-Pereira 2002); openModeller (de Souza Muñoz *et al.*  
101 2011)), but an ongoing problem is that many researchers treat them like “black boxes”,  
102 even though documentation exists in the literature (Joppa *et al.* 2013). In addition to the  
103 shortfalls mentioned above, these GUIs lack adequate guidance within the software.  
104 Further, relying on CLIs for distribution modeling can be challenging even for specialists  
105 because it involves a combination of map inspection, spatial analysis, and statistical  
106 modeling.

107

108 Very recently, a number of new distribution modeling applications present exciting  
109 developments in reproducible science and indicate the demand for software that advances  
110 accessibility and collaboration (e.g. De Giovanni *et al.* 2016; Hallgren *et al.* 2016;  
111 Hardisty *et al.* 2016; Naimi *et al.* 2016; Golding *et al.* Accepted). However, those that  
112 highlight customizability and modularity require programming skills and currently lack  
113 integrated guidance on methods, while others that feature user-friendly interfaces and

114 extensive educational resources are less flexible and have fewer opportunities for user  
115 contributions. *Wallace* aims to provide a wide variety of advantages by having an easily  
116 navigable interface featuring advanced and expandable modeling tools, guidance on  
117 theory and methods, and access to the underlying code. Further development of these  
118 innovative applications and cross-collaboration among them—including *Wallace*—would  
119 benefit the field greatly. Clearly no single lab or research group can address all the needs  
120 of biogeography and related fields, and *Wallace* is designed to expand in an agile fashion  
121 as the field advances and new demands arise.

122

### 123 **WALKTHROUGH**

124 We present a brief walkthrough of *Wallace* v1.0, which is divided into a series of  
125 components that each feature one or more modules. Module authors and featured R  
126 packages are documented in each module (Fig. 2). All major modules have associated  
127 unit tests, and these will be standard with module submission going forward.

128

129 **1) Obtain Occurrence Data:** Species occurrence records can be obtained from online  
130 databases or supplied by the user. *Wallace* currently accesses GBIF, VertNet, and BISON,  
131 removes duplicate coordinates, plots localities on a map, and populates a data table.

132 **2) Process Occurrence Data:** The user chooses which localities to include in the  
133 analysis and can address sampling bias by selecting localities on a map or using a spatial-  
134 thinning algorithm (Aiello-Lammens *et al.* 2015).

135 **3) Obtain Environmental Data:** For gridded predictor variables to characterize the  
136 species' response to the environment, *Wallace* currently offers WorldClim bioclimatic  
137 rasters (Hijmans *et al.* 2005) or allows user-input rasters.

138 **4) Process Environmental Data:** The user delineates a study extent to crop the predictor  
139 grids and draw background samples, as required by most presence-only models. *Wallace*  
140 offers four alternatives, with optional buffering: bounding box, minimum convex  
141 polygon, buffers around occurrence points, and user input.

142 **5) Partition Occurrence Data:** To evaluate models, the user chooses among (spatial and  
143 non-spatial) methods to partition occurrence localities into groups for *k*-fold cross-  
144 validation.

145 **6) Build and Evaluate Niche Models:** To examine model complexity, users can fit  
146 multiple models and use evaluation statistics to identify optimal settings (e.g.,  
147 regularization multipliers, feature classes; Hijmans *et al.* 2017; Muscarella *et al.* 2014).

148 **7) Visualize Model Results:** The user can pan around the map to explore suitability  
149 predictions for the study extent, examine response curves for predictor variables, and  
150 view evaluation plots.

151 **8) Project Model:** The user can project models to other areas or time periods. *Wallace*  
152 currently allows future projections based on the estimates of different global circulation  
153 models (Hijmans *et al.* 2005). Critically, users can view the magnitude of environmental  
154 novelty between the study extent and the projected area/time, which can highlight areas  
155 to exercise caution in interpretation (Elith *et al.* 2010).

156 **9) Session Code:** The user can download an R Markdown script that reproduces the  
157 analysis undertaken during the *Wallace* session.

158 **TARGET AUDIENCES**

159 We developed *Wallace* with a wide range of audiences in mind. Graduate students  
160 interested in distribution modeling and coding but who are not yet advanced  
161 programmers should benefit from learning interactively using *Wallace*. Conservation  
162 practitioners and natural resource managers may want to assess data availability and  
163 quality for a study species, learn about methods, run analyses, and share results with  
164 colleagues. Experienced programmers can run models, download the session code, and  
165 customize it to modify or extend the analysis. Those developing new methods may also  
166 want to disseminate their products by contributing new modules to *Wallace*. Lastly,  
167 educators can use *Wallace* to teach interactive lessons about ecology, programming, and  
168 scientific best practices.

169

170 **CONCLUSIONS AND FUTURE DIRECTIONS**

171 *Wallace* demonstrates an innovative, open platform for rapid dissemination of scientific  
172 methods to a broad audience—specifically encouraging plurality in methodology and  
173 ongoing community development. Over the next three years under funding from the U.S.  
174 National Science Foundation, we plan to work closely with a cadry of international  
175 research groups to integrate new modules, both expanding available options within the  
176 existing scope of *Wallace* and broadening the scope of its capabilities. Some plans for the  
177 future include providing more environmental datasets and modeling algorithms,  
178 measuring prediction uncertainty, integrating analyses that use distribution models as  
179 inputs (e.g. measuring biodiversity, conservation planning), and model comparison tools.

180 Above all, our vision for an expandable software like *Wallace* is that users decide what  
181 needs to be added and become contributors themselves.

182 Although *Wallace* is currently focused on distribution modeling, other fields may  
183 benefit from adopting a similar framework for software development. Like distribution  
184 models, many complex analyses can often be broken down into components and  
185 assembled into teachable workflows for disseminating methods to a broad audience.  
186 Furthermore, science advances most quickly when researchers share advancements and  
187 build tools together, which *Wallace* enables. Finally, *Wallace*'s interactive nature  
188 demonstrates an alternative to static manuals, tutorials, or vignettes for presenting new  
189 methods. The next generation of scientific software will benefit from these ideas, which  
190 could lead to more individuals learning, contributing to, and engaged in a dynamic  
191 process of creative collaboration.

192

### 193 **AUTHORS' CONTRIBUTIONS**

194 J.M.K., R.P.A., M.A.-L., B.V., and R.M. conceived of the original genesis of *Wallace*.

195 J.M.K led code development, and B.V., M.A.-L., R.M., and C.M. were co-developers.

196 J.M.K. drafted the manuscript, with major input from R.P.A., and all other coauthors

197 provided revisions. R.P.A. led overall project design and drafted guidance text, on which

198 J.M.K. provided revisions.

199

### 200 **ACKNOWLEDGMENTS**

201 This work was made possible via an award from the Global Biodiversity Information

202 Facility (Finalist, 2015 Ebbe Nielsen Challenge, led by J.M.K.), a CUNY Graduate

203 Center Provost Digital Innovation Grant to J.M.K., and grants from the U.S. National  
204 Science Foundation (DBI-1650241, DBI-1661510, and DEB-1119915 to R.P.A.). B.V.  
205 was supported by a CAPES doctoral grant (Brazil), M.A.-L. by NSF DEB-1046328, and  
206 R.M. by NSF DBI-1401312. We appreciate helpful comments and suggestions from the  
207 Anderson lab and other members of the CUNY Ecology, Evolution, and Behavior group;  
208 Mary E. Blair; and Morgan W. Tingley.

209

#### 210 **DATA ACCESSIBILITY**

211 No data were included in this article. The R package is freely available under the GPL-3  
212 license at Github (development version: <https://github.com/wallaceEcoMod/wallace>) and  
213 CRAN (stable version: <https://cran.r-project.org/package=wallace>). There is also a  
214 project webpage that will be updated with ongoing development  
215 (<https://wallaceecomod.github.io>).

216

217 **REFERENCES CITED**

- 218 1. Aiello-Lammens M.E., Boria R.A., Radosavljevic A., Vilela B. & Anderson R.P.  
219 (2015) spThin: an R package for spatial thinning of species occurrence records for  
220 use in ecological niche models. *Ecography*, 38, 541-545.
- 221 2. Alvarado-Serrano D.F. & Knowles L.L. (2014) Ecological niche models in  
222 phylogeographic studies: applications, advances and precautions. *Molecular*  
223 *Ecology Resources*, 14, 233-248.
- 224 3. Bean W.T., Stafford R. & Brashares J.S. (2012) The effects of small sample size  
225 and sample bias on threshold selection and accuracy assessment of species  
226 distribution models. *Ecography*, 35, 250-258.
- 227 4. Bolker B.M. (2008) *Ecological models and data in R*, Princeton University Press,  
228 Princeton.
- 229 5. Booth T.H., Nix H.A., Busby J.R. & Hutchinson M.F. (2014) BIOCLIM: The first  
230 species distribution modelling package, its early applications and relevance to  
231 most current MaxEnt studies. *Diversity and Distributions*, 20, 1-9.
- 232 6. Chang W., Cheng J., Allaire J.J., Xie Y. & McPherson J. (2017) shiny: Web  
233 Application Framework for R. R package version 1.0.3, [http://cran.r-](http://cran.r-project.org/package=shiny)  
234 [project.org/package=shiny](http://cran.r-project.org/package=shiny).
- 235 7. De Giovanni R., Williams A.R., Ernst V.H., Kulawik R., Fernandez F.Q. &  
236 Hardisty A.R. (2015) ENM Components: a new set of web service-based  
237 workflow components for ecological niche modelling. *Ecography*, 39, 376-383.
- 238 8. Elith J., Kearney M. & Phillips S. (2010) The art of modelling range-shifting  
239 species. *Methods in Ecology and Evolution*, 4, 330-342.

- 240 9. Elith J., Phillips S.J., Hastie T., Dudík M., Chee Y.E. & Yates C.J. (2011) A  
241 statistical explanation of MaxEnt for ecologists. *Diversity and Distributions*, 17,  
242 43-57.
- 243 10. Franklin J. (2010a) *Mapping species distributions: spatial inference and*  
244 *prediction*, Cambridge University Press, New York.
- 245 11. Franklin J. (2010b) Moving beyond static species distribution models in support  
246 of conservation biogeography. *Diversity and Distributions*, 16, 321-330.
- 247 12. Gimenez O., *et al.* (2014) Statistical ecology comes of age. *Biology Letters*, 10,  
248 20140698.
- 249 13. Golding N., August T.A., Lucas T.C.D., Gavaghan D.J., van Loon E.E. &  
250 McNerny G. The zoon R package for reproducible and shareable species  
251 distribution modelling. *Methods in Ecology and Evolution*, Accepted Author  
252 Manuscript, doi:10.1111/2041-210X.12858.
- 253 14. Guisan A. & Thuiller W. (2005) Predicting species distribution: offering more  
254 than simple habitat models. *Ecology Letters*, 8, 993-1009.
- 255 15. Guisan A. & Rahbek C. (2011) SESAM—a new framework integrating  
256 macroecological and species distribution models for predicting spatio-temporal  
257 patterns of species assemblages. *Journal of Biogeography*, 38, 1433-1444.
- 258 16. Hallgren W., *et al.* (2016) The Biodiversity and Climate Change Virtual  
259 Laboratory: Where ecology meets big data. *Environmental Modelling & Software*,  
260 76, 182-186.
- 261 17. Hardisty A.R., *et al.* (2016) BioVeL: a virtual laboratory for data analysis and  
262 modelling in biodiversity science and ecology. *BMC Ecology*, 16, 49.

- 263 18. Hampton S.E., *et al.* (2015) The Tao of open science for ecology. *Ecosphere*, 6,  
264 120.
- 265 19. Hijmans R.J., Cameron S.E., Parra J.L., Jones P.G. & Jarvis A. (2005) Very high  
266 resolution interpolated climate surfaces for global land areas. *International*  
267 *Journal of Climatology*, 25, 1965-1978.
- 268 20. Hijmans R.J., Phillips S., Leathwick J. & Elith J. (2017) dismo: Species  
269 Distribution Modeling. R package version 1.1-4, [http://cran.r-](http://cran.r-project.org/package=dismo)  
270 [project.org/package=dismo](http://cran.r-project.org/package=dismo).
- 271 21. Jiménez-Valverde A., Lobo J.M. & Hortal J. (2008) Not as good as they seem: the  
272 importance of concepts in species distribution modelling. *Diversity and*  
273 *Distributions*, 14, 885-890.
- 274 22. Joppa L.N., McNerny G., Harper R., Salido L., Takeda K., O'Hara K., Gavaghan  
275 D. & Emmott S. (2013) Troubling trends in scientific software use. *Science*, 340,  
276 814-815.
- 277 23. McCormack J.E., Zellmer A.J. & Knowles L.L. (2010) Does niche divergence  
278 accompany allopatric divergence in *Aphelocoma* jays as predicted under  
279 ecological speciation?: Insights from tests with niche models. *Evolution*, 64,  
280 1231-1244.
- 281 24. Merow C., Smith M.J. & Silander J.A. (2013) A practical guide to MaxEnt for  
282 modeling species' distributions: what it does, and why inputs and settings matter.  
283 *Ecography*, 36, 1058-1069.
- 284 25. Mislán K.A.S., Heer J.M. & White E.P. (2016) Elevating the status of code in  
285 ecology. *Trends in Ecology & Evolution*, 31, 4-7.

- 286 26. Muscarella R., Galante P.J., Soley-Guardia M., Boria R.A., Kass J.M., Uriarte M.  
287 & Anderson R.P. (2014) ENMeval: An R package for conducting spatially  
288 independent evaluations and estimating optimal model complexity for Maxent  
289 ecological niche models. *Methods in Ecology and Evolution*, 5, 1198-1205.
- 290 27. Naimi B. & Araújo M.B. (2016) sdm: a reproducible and extensible R platform  
291 for species distribution modelling. *Ecography*, 39, 368-375.
- 292 28. Peterson A.T., Soberón J., Pearson R.G., Anderson R.P., Martinez-Meyer E.,  
293 Nakamura M. & Araújo M.B. (2011) *Ecological niches and geographic*  
294 *distributions*, Princeton University Press.
- 295 29. Phillips S.J., Anderson R.P. & Schapire R.E. (2006). Maximum entropy modeling  
296 of species geographic distributions. *Ecological Modelling*, 190, 231-259.
- 297 30. R Core Team (2017) R: A language and environment for statistical computing. R  
298 Foundation for Statistical Computing. URL: <http://www.R-project.org> [accessed  
299 02 February 2017]
- 300 31. Radosavljevic A. & Anderson R.P. (2014) Making better Maxent models of  
301 species distributions: complexity, overfitting and evaluation. *Journal of*  
302 *Biogeography*, 41, 629-643.
- 303 32. Scachetti-Pereira R. (2002) DesktopGarp: a software package for biodiversity and  
304 ecologic research. United States: The University of Kansas Biodiversity Research  
305 Center. Available online at <http://www.nhm.ku.edu/desktopgarp/>
- 306 33. de Souza Muñoz M.E. *et al.* (2011) openModeller: a generic approach to species'  
307 potential distribution modelling. *GeoInformatica*, 15, 111-135.

- 308 34. VanDerWal J., Shoo L.P., Graham C. & Williams S.E. (2009). Selecting pseudo-  
309 absence data for presence-only distribution modeling: How far should you stray  
310 from what you know? *Ecological Modelling*, 220, 589-594.
- 311 35. Warren D.L & Seifert S.N. (2011) Ecological niche modeling in Maxent: the  
312 importance of model complexity and the performance of model selection criteria.  
313 *Ecological Applications*, 21, 335-342.
- 314 36. Yackulic C.B., Chandler R., Zipkin E.F., Royle J.A., Nichols J.D., Campbell  
315 Grant E.H. & Veran S. (2013) Presence-only modelling using MAXENT: when  
316 can we trust the inferences? *Methods in Ecology and Evolution*, 4, 236-243.
- 317

318 Table 1. Advantages of the *Wallace* framework.

OPEN	<ul style="list-style-type: none"> <li>- code is free and open-source (GNU GPL 3.0)</li> <li>- users can access data from online databases</li> </ul>
EXPANDABLE	<ul style="list-style-type: none"> <li>- modules (discrete methodological options) can be contributed by the community</li> </ul>
FLEXIBLE	<ul style="list-style-type: none"> <li>- multiple options exist for user data uploads and downloads of results</li> </ul>
INTERACTIVE	<ul style="list-style-type: none"> <li>- sessions are participatory and encourage experimentation</li> <li>- a variety of visualizations are provided (maps, tables, figures)</li> </ul>
INSTRUCTIVE	<ul style="list-style-type: none"> <li>- guidance text (theoretical and methodological) is included for all components and modules</li> </ul>
REPRODUCIBLE	<ul style="list-style-type: none"> <li>- an annotated and executable R Markdown file is produced for rerunning analyses, sharing results, providing supplemental information / educational resources</li> </ul>

319