

# Online Research @ Cardiff

This is an Open Access document downloaded from ORCA, Cardiff University's institutional repository: <http://orca.cf.ac.uk/75548/>

This is the author's version of a work that was submitted to / accepted for publication.

Citation for final published version:

De Giovanni, Renato, Williams, Alan R, Hernández Ernst, Vera, Kulawik, Robert, Quevedo Fernandez, Francisco and Hardisty, Alex 2015. ENM Components: a new set of web service-based workflow components for ecological niche modelling. *Ecography -Copenhagen-* 10.1111/ecog.01552 file

Publishers page: <http://onlinelibrary.wiley.com/journal/10.1111/%28...>  
<<http://onlinelibrary.wiley.com/journal/10.1111/%28ISSN%291600-0587>>

Please note:

Changes made as a result of publishing processes such as copy-editing, formatting and page numbers may not be reflected in this version. For the definitive version of this publication, please refer to the published source. You are advised to consult the publisher's version if you wish to cite this paper.

This version is being made available in accordance with publisher policies. See <http://orca.cf.ac.uk/policies.html> for usage policies. Copyright and moral rights for publications made available in ORCA are retained by the copyright holders.



1 ENM Components: a new set of web service-based workflow  
2 components for ecological niche modelling

3 **Renato De Giovanni, Alan R. Williams, Vera Hernández E., Robert Kulawik,**  
4 **Francisco Quevedo Fernandez, Alex R. Hardisty**

5 R. Giovanni ([renato@cria.org.br](mailto:renato@cria.org.br)), Centro de Referência em Informação Ambiental, Av. Dr.  
6 Romeu Tórtima, 388, 13084-791, Campinas, SP, Brazil. – A. R. Williams, School of Computer  
7 Science, University of Manchester, Oxford Road, Manchester, M13 9PL, United Kingdom – V.  
8 H. Ernst and R. Kulawik, Fraunhofer-Gesellschaft zur Förderung der angewandten Forschung  
9 e.V., Postfach 20 07 33, 80007, München, Germany – F. Quevedo and A. Hardisty, Cardiff  
10 School of Computer Science & Informatics, 5, The Parade, Cardiff, CF24 3AA, United Kingdom.

11 **Abstract**

12 Ecological Niche Modelling (ENM) Components are a set of reusable workflow components  
13 specialized for performing ENM tasks within the Taverna workflow management system. Each  
14 component encapsulates specific functionality and can be combined with other components to  
15 facilitate the creation of larger and more complex workflows. One key distinguishing feature of  
16 ENM Components is that most tasks are performed remotely by calling web services,  
17 simplifying software setup and maintenance on the client side and allowing more powerful  
18 computing resources to be exploited. This paper presents the current set of ENM Components  
19 in the context of the Taverna family of tools for creating, publishing and sharing workflows. An  
20 example is included showing how the components can be used in a preliminary investigation of  
21 the effects of mixing different spatial resolutions in ENM experiments.

## 22 **Introduction**

23 By being able to predict and to understand species' distribution under different scenarios,  
24 ecological niche modelling (ENM) recently became one of the most popular techniques in  
25 biodiversity research, with direct impact in the number of published papers (Lobo et al. 2010)  
26 and related tools (see Peterson et al. 2011 for references). Most of the work done in this field  
27 uses the correlative approach (Soberón and Peterson 2005), in which species occurrence  
28 points are combined with environmental data, serving as inputs to a modelling algorithm. The  
29 resulting models can then be projected into different geographical regions under different  
30 environmental scenarios, producing potential distribution maps with a wide range of uses.

31 Although the typical ENM procedure is usually straightforward for a single species with some of  
32 the existing software, many experiments can be quite complex, requiring several steps, usually  
33 mixing different tools. In such cases, a workflow approach through workflow management  
34 systems may offer several benefits. Scientific workflows can specify a sequence of data  
35 retrieval, data manipulation and data storage/publication steps. When a scientific procedure or  
36 protocol is captured as a workflow, this allows the protocol to be easily shareable and re-  
37 runnable. In addition, provenance data of what happened during a workflow run allows for  
38 research to be, within certain limits, reproducible.

39 Considering the two most popular ENM software found by a recent survey (Ahmed et al. 2015),  
40 users seem to be divided between simplicity and flexibility, as if these two features would be  
41 irreconcilable in the same software. That is, if users are looking for an easy to use interface with  
42 a short learning curve, they must live with inflexible point-and-click software, whereas if they  
43 wish flexibility, they must develop programming skills to use syntax driven software. None of the  
44 tools found by the survey are based on workflow management systems, which actually have the  
45 potential to provide both a simple and flexible interface. The creation of scientific workflows is  
46 commonly carried out within a graphical user interface which may be desktop based, for

47 example Taverna Workbench<sup>1</sup> (Wolstencroft et al. 2013) and Kepler (Altintas et al. 2004); or  
48 browser based, for example Taverna Online<sup>2</sup> and Galaxy (Giardine et al. 2005). Such interfaces  
49 allow users to visually build custom workflows, usually by means of adding boxes on a panel  
50 (each box representing a task) and connecting them through input/output parameters. This  
51 intuitive way to design and control personalised workflows is one of the main reasons for  
52 scientific workflows to be currently used in a large number of disparate domains, for example  
53 bioinformatics, astronomy and preservation of digital resources.

54 Most workflow systems allow different types of steps to be included within a workflow, such as  
55 running user-defined scripts, interacting with the user to display or get data, and calling external  
56 programs locally or remotely. In this last case, workflows may perform tasks by interacting with  
57 web services. Web services are software applications supporting dynamic interactions with  
58 other programs over the Internet through open standards. Using web services inside workflows  
59 may bring up issues related to the need of having an Internet connection and to the reliability  
60 and limitations of third-party service providers. However, web services also offer considerable  
61 advantages in terms of minimising the need for software installation and maintenance on the  
62 client side. There can also be more powerful computational resources behind web services,  
63 allowing workflows to outsource part of the processing requirements and not be strictly  
64 constrained by a desktop environment.

65 The Taverna suite of tools is a workflow management system allowing the creation, editing,  
66 sharing and running of workflows. Taverna workflows may be created and edited within the  
67 desktop Taverna Workbench or using the web-based Taverna Online. Workflows may be run: 1)  
68 directly within Taverna Workbench, 2) locally by the Taverna Command Line Tool or 3) remotely  
69 on the Taverna Server, which allows multiple simultaneous runs with secure user separation  
70 and offers a web service interface that can be remotely invoked by other programs. Additionally,

---

<sup>1</sup> <http://www.taverna.org.uk>

<sup>2</sup> <http://onlinehpc.com>

71 the running of a Taverna workflow can be included within a web application by using a self-  
72 contained software package called Taverna Player, similar to the way that videos are currently  
73 embedded within web applications. Taverna Player handles the marshalling of input data to and  
74 results from runs on a Taverna Server, also handling interaction requests from workflow runs.  
75 Taverna Player can be included in diverse web applications, such as IPython Notebook (Pérez  
76 & Granger 2007), Scratchpad (Smith et al. 2011) and web portals. Finally, any Taverna  
77 workflow can also be easily shared in the myExperiment platform<sup>3</sup>.

78 Some of the recent developments in Taverna were carried out as part of the Biodiversity Virtual  
79 e-Laboratory (BioVeL) project<sup>4</sup>. BioVeL placed particular emphasis on setting up a robust web  
80 service infrastructure upon which scientific workflows can be built. This effort involved improving  
81 existing web services and creating new ones when necessary. All web services being used by  
82 BioVeL are registered in the Biodiversity Sciences Web Service Catalogue<sup>5</sup>, including service  
83 endpoint, documentation, and monitoring information.

## 84 **Workflows & ENM**

85 Historically, ENM has been among case studies in many projects focused on scientific  
86 workflows. In 2004 the Biodiversity World project used the Triana workflow management system  
87 (Taylor et al. 2003) to build ENM workflows (Pahwa et al. 2006). Almost in parallel, the Science  
88 Environment for Ecological Knowledge (SEEK) project also created ENM workflows  
89 (Pennington et al. 2007), this time using the Kepler system. More recently, the ENM workflow  
90 approach was revived with the SAHM module in VisTrails (Morissette et al. 2013) and with the  
91 BioVeL project, where ENM is one of the major research areas (see Leidenberger et al. 2014 for  
92 an example).

---

<sup>3</sup> <http://myexperiment.org>

<sup>4</sup> <http://biovel.eu>

<sup>5</sup> <http://biodiversitycatalogue.org>

93 Even with all these initiatives, workflow management systems are still seen as a rather  
94 challenging environment for most researchers, traditionally requiring significant programming  
95 expertise to perform any different task that is typically needed when creating a custom workflow.  
96 Moreover, without sufficient specific analytical functions and features needed by ecologists and  
97 biodiversity researchers, the familiarisation effort required from researchers to start using  
98 workflow tools has not yet been perceived as sufficiently worthwhile. To overcome these  
99 challenges, one of the approaches explored at BioVeL has been to create families of workflow  
100 components specialized in common tasks for a certain area, such as ENM or phylogenetics.  
101 Each component is a sub-workflow representing a task-unit encapsulating implementation  
102 details. Components offer a high-level interface, allowing them to be more easily used and  
103 combined to create larger workflows.

104 BioVeL also created a web portal<sup>6</sup> where users can upload workflows or reuse workflows  
105 uploaded by other users. The portal allows users to start multiple workflow runs and retrieve  
106 results later, without needing an active Internet connection during the workflow run when there  
107 is no interaction involved. There are no additional requirements for a user to run a workflow  
108 through the portal except having an Internet browser.

109 Another major concern in BioVeL was to assure sustainability of assets beyond the project  
110 lifetime – especially considering that most of its workflows are strongly based on web services  
111 provided by different institutions. BioVeL's strategy to maintain a stable and persistent e-  
112 Infrastructure largely depends on institutional commitment, where each individual organisation  
113 takes responsibility to sustain various pieces of the e-Infrastructure as part of its core business.  
114 A typical example is the ENM service provided by the Reference Center on Environmental  
115 Information (CRIA), which is currently used by ENM Components. The service has been running  
116 for many years at CRIA, well before the BioVeL project started, and will continue to run, as it is  
117 considered an important asset fully aligned to the institutional mission. Still regarding

---

<sup>6</sup> <http://portal.biovel.eu>

118 sustainability, BioVeL satisfies two pre-requisites pointed out by Henfridsson and Bygstad  
119 (2013) as being important factors for the adoption, spreading and evolution of a digital e-  
120 Infrastrucure: 1) loosely-coupled, service-oriented architecture and 2) decentralised  
121 management. All these factors contribute to the availability and improvement of ENM  
122 Components over time.

## 123 **ENM Components**

124 The ENM Components were created with the Taverna workflow management system as part of  
125 BioVeL to simplify the existing ENM workflows produced by the project and to facilitate the  
126 creation of new workflows. Since Taverna components are special workflows themselves, they  
127 enjoy the same benefits of the Taverna suite: they can be designed and run using the same  
128 tools, they can be reused by other workflows and even shared in myExperiment, where ENM  
129 Components are all publicly available under a specific pack with the same name<sup>7</sup> (note: to use  
130 them it is not necessary to manually download the pack, as Taverna Workbench can  
131 dynamically interact with myExperiment to fetch remote components).

132 A main aspect of providing reusable components is to document how they can be used. In this  
133 respect, each ENM component has a short description of its functionality and of each  
134 input/output parameter (also called ports). Being workflows, all components can be opened with  
135 Taverna workbench and run independently (all ports provide example values that can be used  
136 for testing). Using ENM Components to build new workflows within the workbench is only a  
137 matter of dragging the desired component from the service panel into the workflow being  
138 designed. To facilitate the connecting of different components, most ports with equal interfaces  
139 (same parameters and data types) are assigned the same name (fig. 1 shows how the main  
140 components can be combined). More information about how to use ENM Components can be  
141 found in the corresponding pack description in myExperiment. All ENM workflows developed in

---

<sup>7</sup> <http://www.myexperiment.org/packs/563>

142 BioVeL are based on ENM Components, providing many examples of their usage (see the  
143 generic “Ecological niche modelling workflow”<sup>8</sup> and the “Bioclim workflow”<sup>9</sup>).

144 The web service currently used by ENM Components was developed on top of openModeller  
145 (Muñoz et al. 2011). OpenModeller is a toolbox mainly comprised by an ENM framework with a  
146 comprehensive list of functions that can be called by other programs. The framework has many  
147 algorithms available and makes use of other software libraries to handle different data formats  
148 and spatial reference systems. The openModeller toolbox also contains a set of command-line  
149 tools and the web service itself, both making use of the framework and sharing most data  
150 structures for input/output parameters. Since ENM Components are strongly based on the  
151 openModeller web service, sometimes it may be necessary to refer to the web service  
152 documentation<sup>10</sup> when designing new workflows. For example, many ports of the ENM  
153 Components return or expect data according to openModeller serialization rules. The three-  
154 tiered structure currently used by ENM Components (component/web service/server software)  
155 actually allows for alternative implementations in the future, provided the same input/output  
156 ports remain the same for each component. For instance, a different web service  
157 implementation could be used (not necessarily associated with openModeller tools), or even all  
158 web service calls could be replaced with interactions to locally installed software. At the  
159 moment, the implementation of ENM Components takes advantage of all algorithms available in  
160 openModeller and of its capabilities to handle different data formats and spatial reference  
161 systems, interacting with a remote web service provided by CRIA.

162 Using remote web services in ENM tasks brings a few changes in the way researchers are used  
163 to working with traditional stand-alone tools. For example, the service needs to be queried to  
164 know which algorithms can be used. Over time, new or enhanced algorithms may become  
165 available on the service being called (information about the currently available algorithms can

---

<sup>8</sup> <http://www.myexperiment.org/workflows/3355>

<sup>9</sup> <http://www.myexperiment.org/workflows/3725>

<sup>10</sup> [http://openmodeller.sf.net/web\\_service\\_2.html](http://openmodeller.sf.net/web_service_2.html)



166 also be found in the openModeller web site<sup>11</sup>). Frequently used environmental layers are  
167 available on the server for convenience, and again the service can be queried to return this  
168 information. Alternatively, additional layers or masks can be provided to the service, as the  
169 modelling engine can access other geospatial web services such as WCS<sup>12</sup> or remote files  
170 available over the web. In this case, layers need to be uploaded somewhere, for example a  
171 BioSTIF<sup>13</sup> server. BioSTIF provides a set of standardized services for spatial data visualization,  
172 transformation and storage. Some of the ENM Components rely on BioSTIF to visualize points  
173 and maps.

#### 174 **Example: the effect of mixing different spatial resolutions**

175 During the BioVeL project, one of the case studies faced a common situation in ENM:  
176 environmental layers came from different sources in different spatial resolutions, i.e., having  
177 different cell sizes (see Leidenberger et al. 2014 for more details). Although the sensitivity of  
178 models to spatial resolution has already been studied before (Guisan et al. 2007), we could not  
179 find specific references about mixing layers with different resolutions. Probably the main reason  
180 is that most of the existing ENM software actually forces researchers to use layers with exactly  
181 the same resolution, spatial reference system and extent – even when differences are  
182 negligible. Since openModeller does not have this constraint – and consequently also the  
183 service used by the ENM Components – users are left with the decision about what to do when  
184 there are such differences between layers.

185 The usual practice when environmental datasets come in different resolutions is to previously  
186 downscale the low resolution layers by subdividing their cells, or upscale the high resolution  
187 ones by coarsening their cell size. It is also important to note that two main factors should be  
188 taken into account when dealing with spatial resolution in an ENM experiment: 1) the resolution

---

<sup>11</sup> <http://openmodeller.sf.net/documentation.html>

<sup>12</sup> <http://www.opengeospatial.org/standards/wcs>

<sup>13</sup> <http://www.biodiversitycatalogue.org/services/7>

189 of the biological phenomenon being studied, since each species may respond to environmental  
190 signals at different scales (Peterson et. al 2011) and 2) the spatial uncertainty of the occurrence  
191 points being used. Ideally, this uncertainty should not be greater than the environmental cell  
192 size, otherwise models will be generated with mistakenly precise environmental data, which  
193 tends to degrade model performance (Graham et al. 2008). Thus, when both factors are  
194 compatible with the finest environmental resolution at hand, which approach – downscaling or  
195 upscaling environmental layers – produces better models? In this example, different features of  
196 the ENM Components and the workflow approach are demonstrated, showing a possible way to  
197 investigate the subject. In particular, we demonstrate the flexibility and modularity of ENM  
198 Components combining them in a workflow that contains user interaction, loop, custom code  
199 and more than one tool. The workflow also explores some of the capabilities of openModeller,  
200 such as generating virtual niches and handling environmental layers in different resolutions.

201 OpenModeller handles differences in spatial resolution and reference systems by treating each  
202 layer independently and simply fetching the corresponding environmental data at each point  
203 (presence, absence or background). Therefore, mixing layers with different spatial resolutions in  
204 openModeller is essentially equivalent to downscaling the low resolution layers with the nearest-  
205 neighbor method, which retains the same original cell value in the new smaller cells. The only  
206 difference with other ENM software is that there is no raster downscaling pre-processing step  
207 inside or outside openModeller – it uses the original layers without modifications. In this  
208 example, we simulate the situation of having environmental layers in different spatial resolutions  
209 and compare the results of models generated with the original layers (mixed resolutions,  
210 equivalent to downscaling the low resolution layers) with models generated after upscaling the  
211 high resolution layers. The workflow created can be summarized in eleven steps (fig. 2), with  
212 the first step involving user interaction to choose the environmental layers and study area  
213 (mask), followed by a loop containing most steps, including virtual niche generation, point

214 sampling, model creation and model testing. A final step after the loop compares the results and  
215 generates a graph using another tool.

216 The workflow can be downloaded from myExperiment<sup>14</sup> and requires: Taverna 2.5; R (R  
217 Development Core Team 2008) with the Rserve package installed and running as a localhost  
218 service in the default port (6311); an active Internet connection so that the workflow can  
219 communicate with the external ENM service currently hosted at CRIA<sup>15</sup>, and a web browser to  
220 handle user interactions. The R version used was 3.1.2 and the Web browser was Firefox  
221 34.0.5. When running this workflow using the Taverna Workbench with the default values  
222 (10000 background points and 30 iterations), it is highly recommended to disable provenance  
223 capture and in-memory storage in the system preferences. The workflow run takes about an  
224 hour to complete with the current resources on the web service, but it may take longer  
225 depending on connection and service load. A simplified version of the workflow with a single  
226 iteration and including model projections is also available<sup>16</sup>.

227 The basic idea of the workflow is to compare models generated with mixed resolution layers  
228 (downscaling scenario) with models generated only with low resolution layers (upscaling  
229 scenario), testing them against the same set of points extracted from a virtual species niche.  
230 The workflow initially retrieves all available layers on the server and asks the user to choose a  
231 set of environmental layers and then a mask delimiting the study area. This initial step is  
232 performed by a nested workflow labeled “choose layers and mask”, containing only a few  
233 interconnected ENM Components and constant values used as input parameters. Each kind of  
234 workflow element in Taverna has a different background color and any workflow element can be  
235 renamed. Components are displayed with a pink background, and most ENM Components used  
236 by this workflow were renamed to better indicate their purpose (original names can always be  
237 found in the details of each component). There are currently a few mask options offered by the

---

<sup>14</sup> <http://www.myexperiment.org/workflows/4535>

<sup>15</sup> <http://modeller.cria.org.br/ws2/om>

<sup>16</sup> <http://www.myexperiment.org/workflows/4536>

238 ENM service, all of them based on political boundaries, which does not affect an arbitrary mask  
239 choice for this study. For simplicity, it can even be assumed for any chosen mask that the whole  
240 area has been historically accessible to the virtual species that will be created in one of the next  
241 steps, so that presence points can be undoubtedly interpreted as being suitable for the species,  
242 and absence points unsuitable. After choosing a mask, the user is then asked to select a set of  
243 high resolution environmental layers, and in the next step to pick the corresponding low  
244 resolution ones. The choice of environmental layers is also arbitrary, and we can also assume  
245 that the chosen layers are the main variables that determine the virtual species' niche. For the  
246 purpose of this experiment, the only constraint when choosing layers is to select variables that  
247 are available at least in two different resolutions. Worldclim bioclimatic variables (Hijmans et al.  
248 2005) are available on the ENM service in 30 arc-seconds and 10 arc-minutes resolutions,  
249 making them a convenient choice for the demonstration. Additionally, Worldclim layers were  
250 originally produced in 30 arc-seconds, with all other low resolution versions – including the 10  
251 arc-minutes one – being obtained by upscaling (Hijmans et al. 2005). At the end of this initial  
252 step, the workflow has two sets of environmental layers with exactly the same variables, each  
253 set with a different spatial resolution.

254 Next, the workflow uses another ENM Component to randomly sample 10000 background  
255 points over the whole study area. At this stage, two other elements are used to demonstrate  
256 how to include custom code in a workflow. Depicted in brown background, they are known as  
257 Java BeanShells<sup>17</sup> in Taverna. One of them (“merge all layers”) concatenates the identifiers of  
258 all selected layers in a single string list before sampling background points, to ensure that all  
259 sampled points have valid values across all different layers and resolutions. The other one (“for  
260 loop triggering”) simply creates a list with the same size of the “replicates” workflow parameter,  
261 since workflow loops can be activated by lists. Although custom code may require programming  
262 skills, BeanShells can easily be transformed into new workflow components if necessary, and

---

<sup>17</sup> <https://jcp.org/en/jsr/detail?id=274>

263 then stored in specific components families to be used by other users. A few additional  
264 BeanShell examples can be found in other parts of the workflow.

265 Next, we use a workflow loop to repeat the same steps a specified number of times. These  
266 steps are inside the “create and test models” nested workflow, where a virtual niche is  
267 generated, training and testing points are sampled based on the virtual species distribution, and  
268 finally the two models for each set of layers are generated and tested. In the first part, the ENM  
269 Component for sampling points is used again to sample a single point to be passed as a  
270 parameter to the Virtual Niche algorithm in openModeller. This algorithm assumes that the  
271 corresponding environmental values for the point are the optimal conditions for the virtual  
272 species, randomly defining standard deviations for each variable to create a continuous niche  
273 across the study area. This is all performed with the high resolution environmental layers,  
274 producing a high resolution niche to be considered the truth for the virtual species. The  
275 corresponding niche is then evaluated over all background points to get the niche values, which  
276 are ordered and split based on a random threshold separating suitable from unsuitable  
277 conditions, ensuring a random arbitrary prevalence between 0.1 and 0.7. These two groups of  
278 points (suitable/unsuitable) are used to randomly sample presence points for model creation (a  
279 number between 30 and 100) and 100 points for independent model testing (50 presences and  
280 50 absences). Finally, the workflow creates two models using one of the most popular ENM  
281 algorithms also available in openModeller: Maxent (Phillips et al. 2006). The first model is  
282 created with the corresponding low resolution environmental layers (upscaling scenario) and the  
283 other with a random balanced mix of high and low resolution layers (downscaling scenario).  
284 These models are tested with the same testing points by measuring the area under the  
285 Receiver Operating Characteristic curve (AUC) – a threshold-independent test suitable for  
286 algorithms that produce a continuous (non binary) output such as Maxent. AUC values range  
287 from 0 to 1, where 1 indicates perfect discrimination between the presence and absence points  
288 being tested, 0.5 indicates a predictive discrimination equivalent to a random guess, and values

289 below 0.5 indicate discrimination worse than random. All steps from virtual niche generation  
290 until model tests are repeated 30 times in the workflow to generate enough variation in the  
291 virtual niche, training points, testing points and resolution mix with the selected layers. In the last  
292 part of the workflow, results are compared using an R script which also produces a graph  
293 plotting side by side AUC values for each set of layers in each iteration. This way, the example  
294 also demonstrates how to use different tools in different parts of the same workflow. The  
295 probability (p-values) of getting a better model when mixing resolutions (downscaling scenario)  
296 instead of using only low resolution layers (upscaling scenario) is estimated as the percentage  
297 of times that the former AUC is greater than the later one. This is a two-tailed test also used by  
298 Elith et al. (2006) to compare the performance of different algorithms. A value close to 1 means  
299 that mixing resolutions produces better models than using only low resolution layers, and vice-  
300 versa for a value close to 0.

301 In the first workflow run, we used Mexico as the mask and WorldClim bio2 (mean diurnal range),  
302 bio5 (maximum temperature of warmest month), bio6 (minimum temperature of coldest month),  
303 bio12 (annual precipitation) and bio14 (precipitation of driest month) as the environmental  
304 variables. Most models using mixed layer resolutions produced better AUCs, although the  
305 differences were small (fig. 3) and the result was not significant ( $p=0.73$ ). We also used the  
306 simplified version of the workflow with the same parameters to project models, illustrating a  
307 virtual species distribution (fig. 4) and its corresponding projected model with mixed resolutions  
308 (fig. 5). Back to the complete workflow, an identical pattern was found in a subsequent run with  
309 different parameters: Finland as the mask and bio2, 3, 4, 6, 13 and 14 as the environmental  
310 variables ( $p=0.73$ ). A third run using India as the mask and bio1, 4, 11, 15 and 16 as the  
311 environmental layers pointed to the same direction, but with less intensity ( $p=0.53$ ).

312 Since the main purpose of the example was to demonstrate the use of ENM Components, we  
313 tried not to add more complexity to the workflow. For a more extensive investigation, future  
314 versions of the workflow could for example include automatic variation of mask, number of

315 layers and proportion of mixed layers, also including more spatial resolutions. An additional step  
316 to produce biased training points could produce a wider and more realistic range of AUCs.  
317 Other modelling algorithms could be tested as well.

318 Even being just a preliminary investigation, the example shows how the ENM Components can  
319 be combined to produce unique scientific workflows. Additionally, the workflow also shows how  
320 to include other tools into the same workflow, such as the currently ubiquitous R, and how to  
321 include custom code, which can be transformed into new components whenever necessary.  
322 Another possibility for new workflows is to combine components from different areas, such as  
323 the phylogenetics components also created during the BioVeL project, or to benefit from other  
324 Taverna-related tools, such as the workflow parameter optimization plug-in that can be used  
325 with ENM (Holl et al. 2013). There are still many practical uses and research opportunities in  
326 ENM that can be explored, and we hope that ENM Components can provide a flexible and  
327 powerful alternative for future works in this area.

328 *Acknowledgements* – This work is part of the BioVeL project with funding from the European  
329 Union's Seventh Framework Programme for research, technological development and  
330 demonstration under grant agreement no. 283359.

## 331 **References**

- 332 Ahmed, S.E. et al. 2015. Scientists and software – surveying the species distribution modelling  
333 community. – *Divers Distrib*, 21(3): 258–267.
- 334 Altintas, I. et al. 2004. Kepler: an extensible system for design and execution of scientific  
335 workflows. – in: *Proc 16th Int Conf Scientific and Statistical Database Manage*, pp.423–424.
- 336 Elith, J. et al. 2006. Novel methods improve prediction of species' distributions from occurrence  
337 data. – *Ecography*, 29(2): 129–151.
- 338 Giardine, B. et al. 2005. Galaxy: a platform for interactive large-scale genome analysis. –  
339 *Genome Res*, 15(10): 1451–1455.

340 Graham, C. H. et al. 2008. The influence of spatial errors in species occurrence data used in  
341 distribution models. – *J Appl Ecol*, 45(1): 239–247.

342 Guisan, A. et al. 2007. Sensitivity of predictive species distribution models to change in grain  
343 size. – *Divers Distrib*, 13(3): 332–340.

344 Henfridsson, O. and Bygstad, B. 2013. The generative mechanisms of digital infrastructure  
345 evolution. – *Mis Quarterly* 37(3): 907–931.

346 Hijmans, R. J. et al. 2005. Very high resolution interpolated climate surfaces for global land  
347 areas. – *Int J Climatol*, 25: 1965–1978.

348 Holl, S. et al. 2013. On specifying and sharing scientific workflow optimization results using  
349 research objects. – in: *Proc 8th Workshop on Workflows in Support of Large-Scale Science*,  
350 New York, ACM Press, pp.28–37.

351 Leidenberger, S. et al. 2014. Mapping present and future predicted distribution patterns for a  
352 meso-grazer guild in the Baltic Sea. – *J Biogeogr*, advance online publication, doi:  
353 10.1111/jbi.12395

354 Lobo, J. M. et al. 2010. The uncertain nature of absences and their importance in species  
355 distribution modelling. – *Ecography*, 33:103–114.

356 Morissette, J. T. et al. 2013. VisTrails SAHM: visualization and workflow management for species  
357 habitat modeling. – *Ecography* 36: 129–135.

358 Muñoz et al. 2011. openModeller: a generic approach to species' potential distribution  
359 modelling. – *Geoinformatica*, 15: 111–135.

360 Pahwa, J. S. et al. 2006. Biodiversity World: A Problem-Solving Environment for Analysing  
361 Biodiversity Patterns. – in: *Proc 6th IEEE Int Symp Cluster Computing and the Grid (CCGRID)*,  
362 Singapore, pp. 201–208.



363 Pennington, D. D. et al. 2007. Ecological Niche Modelling Using the Kepler Workflow System. –  
364 in: Taylor, I.J., Deelman, E., Gannon, D.B., Shields, M. (Eds.), *Workflows for e-Science*,  
365 Springer, Berlin, pp.91–108.

366 Pérez, F. and Granger, B. E. 2007. IPython: a system for interactive scientific computing. –  
367 *Comput Sci Eng*, 9(3): 21–29.

368 Peterson, A. T. et al. 2011. *Ecological Niches and Geographical Distributions*. – Princeton  
369 University Press.

370 Phillips, S. J. et al. 2006. Maximum entropy modelling of species geographic distributions. –  
371 *Ecol Model*, 190: 231–259.

372 R Development Core Team. 2008. *R: A language and environment for statistical computing*. – R  
373 Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0, URL [http://www.R-](http://www.R-project.org)  
374 [project.org](http://www.R-project.org)

375 Smith, V. S. et al. 2011. Scratchpads 2.0: a Virtual Research Environment supporting scholarly  
376 collaboration, communication and data publication in biodiversity science. – *Zookeys*, 150: 53–  
377 70.

378 Soberón, J. and Peterson A. T. 2005. Interpretation of models of fundamental ecological niches  
379 and species' distributional areas. – *Biodiversity Informatics* 2: 1–10.

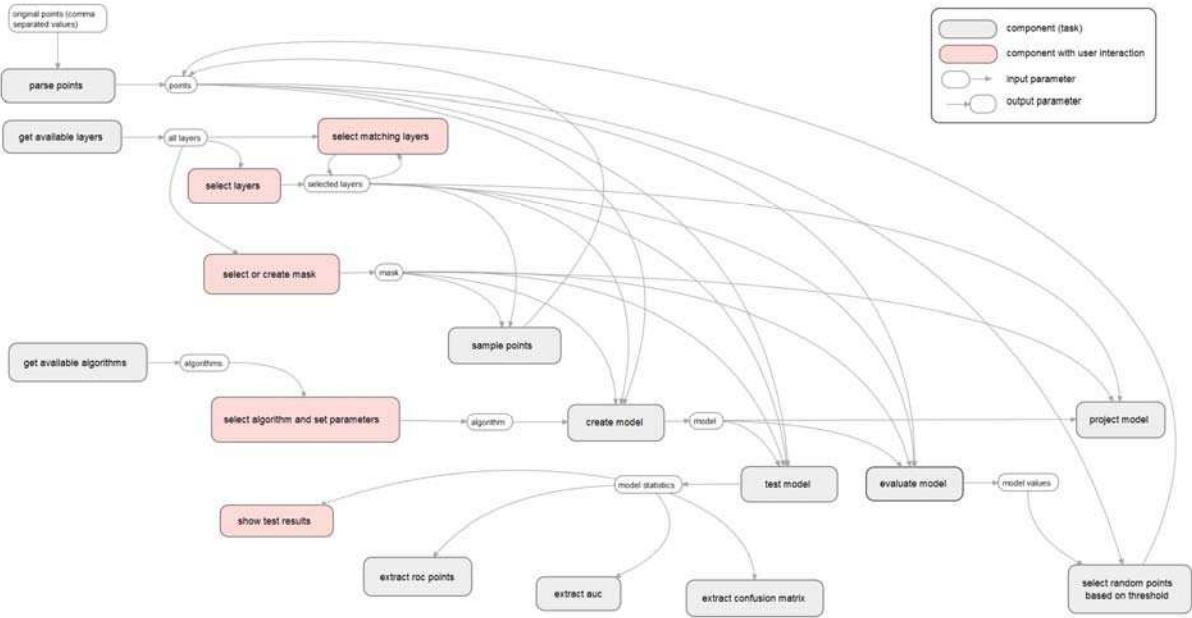
380 Taylor, I. et al. 2003. Triana Applications within Grid Computing and Peer to Peer Environments.  
381 – *J Grid Comput* 1(2): 199–217.

382 Wolstencroft, K. et al. 2013. The Taverna workflow suite: designing and executing workflows of  
383 Web Services on the desktop, web or in the cloud. – *Nucleic Acids Res* 41(1): 557–561.

384

385 **Figures**

386 **Figure 1:**

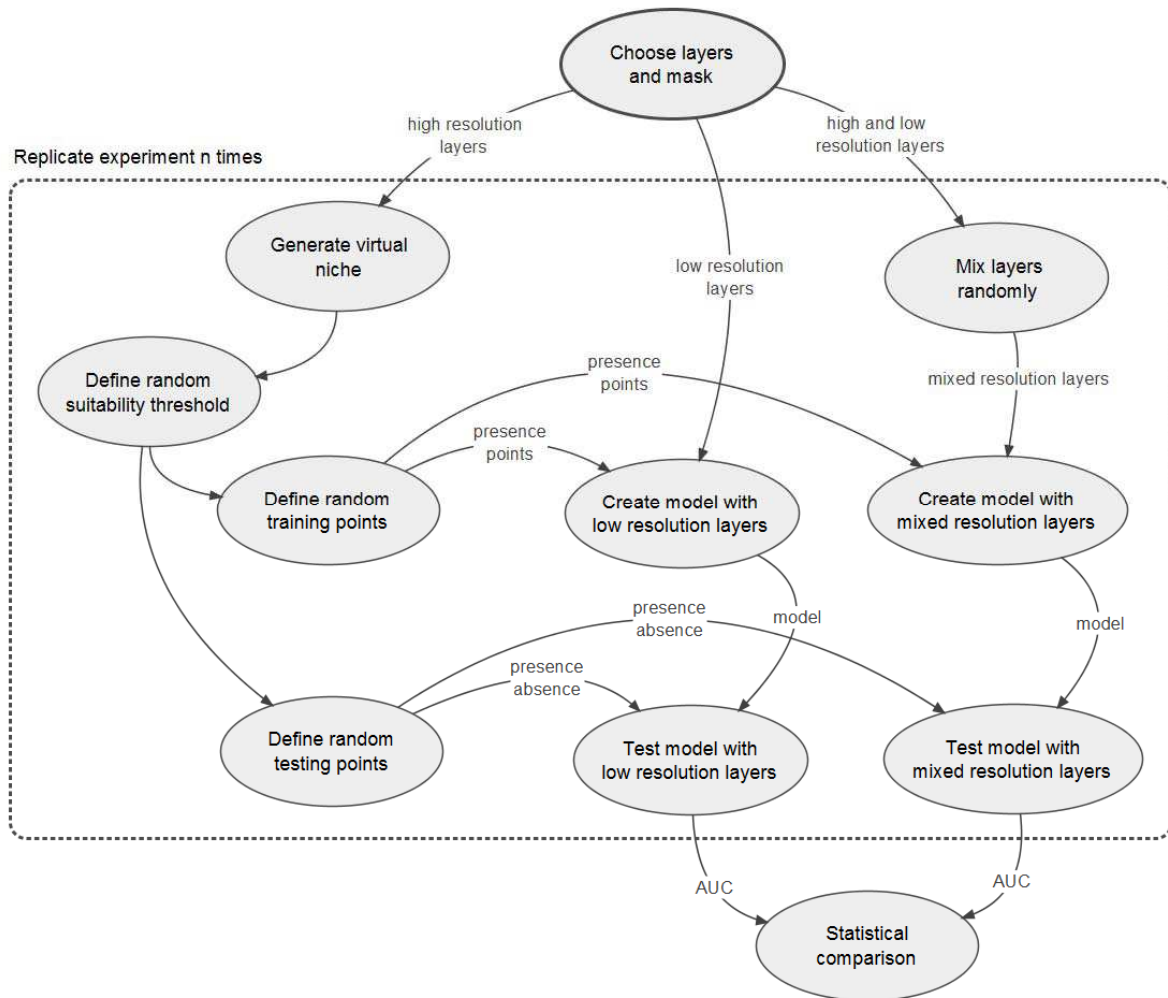


387

388

389

390 Figure 2:

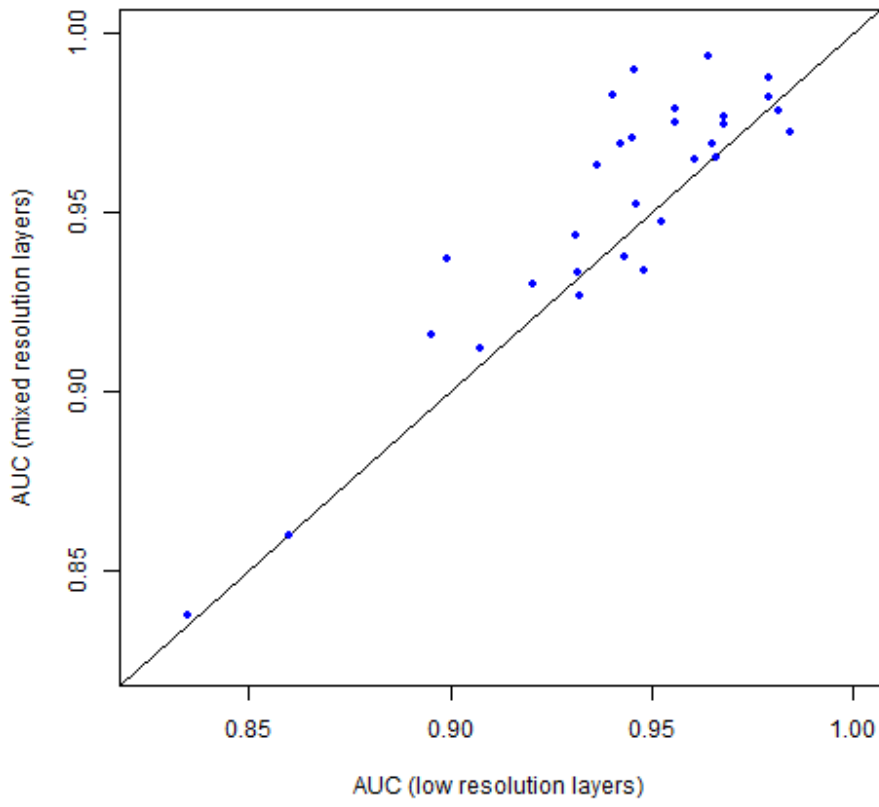


391

392

393

394 Figure 3:

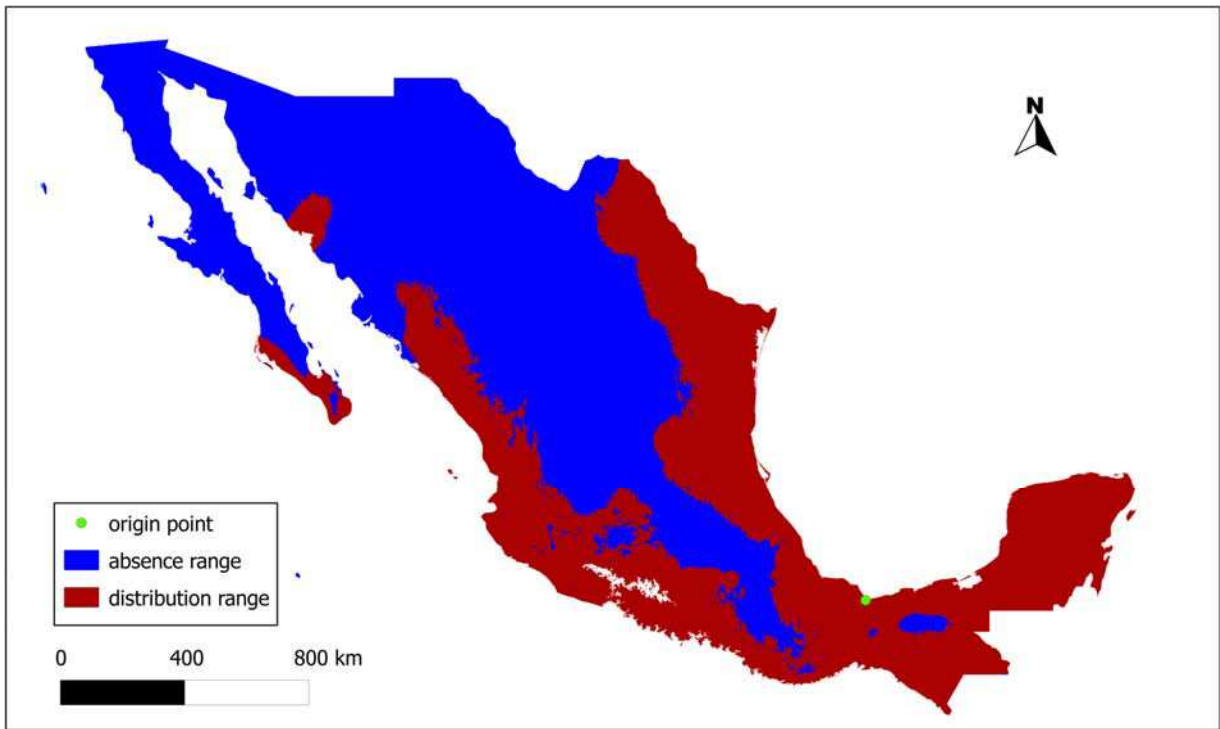


395

396

397

398 Figure 4:

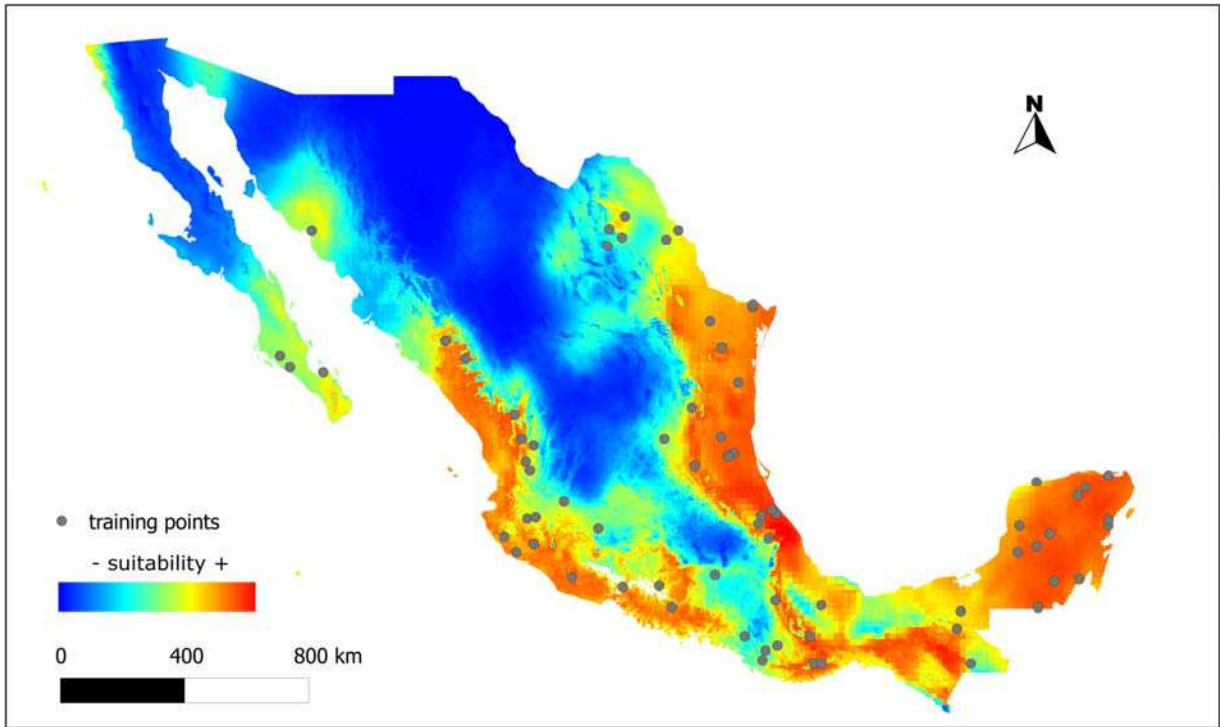


399

400

401

402 Figure 5:



403

404

405