

Computational resources for simulating under a spatial coalescent model across heterogeneous landscapes and testing hypotheses about the geography of genetic variation: QUETZAL-EGGS, -CRUMBS, -NEST and DECRYPT

Arnaud Becheler^{1*} | L. Lacey Knowles^{1†}

¹Ecology and Evolutionary Biology
Department, University of Michigan,
48109, MI, U.S.

Correspondence

Lacey L. Knowles, Ecology and Evolutionary
Biology Department, University of
Michigan, 48109, MI, U.S.
Email: knowlesl@umich.edu

Funding information

This study was funded by NSF [DEB
16-55607 to Lacey L. Knowles

Spatially explicit coalescent models in which the underlying demographic parameters are informed by the environment (either past, present, or temporally and spatially changing environments) provide a framework for hypothesis testing that incorporates geographic information about genetically sampled individuals. This general approach - Integrated Distributional, Demographic and Coalescent (iDDC) modelling - can be used to explain how heterogeneous, dynamic landscapes shape the history and genetic patterns of a species. However, iDDC approaches involve long and complex tasks that often require custom-fit simulators, some coding expertise, and extensive computing resources. Here we introduce several resources that offer improved speed and generality, as well as expand the feasible parameter space for conducting iDDC analyses compared to other software applications. Specifically, QUETZAL-EGGS are C++ iDDC simulators; QUETZAL-CRUMBS is a complementary set of

* A.B. designed the models, software and computational framework.

† L.L.K. supervised the conception, funding and findings of this work.

*† All authors provided critical feedback and helped shape the research, resources and analysis. All authors contributed to the final manuscript.

Python tools for simulating on specific landscapes and conducting Approximate Bayesian Computation (ABC) analyses (e.g., prior sampling, geospatial operations, ENM/SDM, visualization); DECRYPT is a framework for automated, biology informed robustness analysis of the multispecies coalescent model. All these tools and their dependencies for local use or remote computations are made readily available in a Docker container package called QUETZAL-NEST.

KEYWORDS

environmental niche modeling, coalescence, biogeography, software, simulation, landscape

1 | INTRODUCTION

Integrating distributional, demographic and coalescence models (*iDDC modeling*, He et al., 2013, see also the excellent review by Larsson et al. 2021) is a powerful tool to explore how spatial and temporal landscape heterogeneity shapes the genetic diversity of modern populations (e.g., Knowles and Alvarado-Serrano, 2010; Brown and Knowles, 2012; Pan et al., 2020). In this category of modelling approaches, the landscape is discretized into a very large number of demes (e.g., >1000). The demographic history (i.e., the number of individuals in each deme and the number of migrants across demes) is simulated as a function of the environmental variation over many generations (i.e., thousands, to tens of thousands, of generations). Then conditionally on this historical demographic processes, a coalescence process tracks the ancestry backward in time of genetically sampled individuals across a landscape.

When coupled with simulation-based inference methods like Approximate Bayesian Computation, ABC (Beaumont et al., 2002; Estoup et al., 2010), these *iDDC* models have the interesting property to generate complex geographic distributions of genetic variation while maintaining a reduced number of parameters (because parameters describe landscape-wide processes as a function of the underlying environment, rather than excessive parameterization of each individual deme). With less than a dozen parameters, the procedure is nevertheless flexible enough to represent reasonably complex processes (e.g., shifting species distributions, varied rates of migration across a landscape, population growth, and geographic barriers that vary in their attenuation of gene flow).

However, overall contributions of *iDDC* modeling have been rather limited. Few researchers apply the approach despite its intriguing potential for hypothesis testing using biologically informed expectations, and even though *iDDC* modeling addresses questions that could not otherwise be addressed with generic models that are not spatially explicit (e.g., the contribution of contemporary versus historical landscapes to genetic structure; He et al., 2013) recolonization of river routes following deglaciation (Neuenschwander et al., 2008); the geographic position of refugial populations (Bemmels et al., 2019); the facilitative versus competitive effects of co-distributed species on colonization of landscapes (Ortego and Knowles, 2020). We argue that the limited traction of *iDDC* modeling reflects technical and practical challenges of *iDDC* modeling itself. To increase the accessibility of *iDDC* modeling to a broad audience, we have developed a set of software tools that solve some of the methodological hurdles associated with ease of application and computational constraints. Rather than presenting a biological application of these resources in this article, we instead direct readers to an example repository with a full data analysis that is accompanied by detailed documentation

56 of the applied software and computational resources https://github.com/Becheler/quetzal_on_OSG.

57 2 | QUETZAL-EGGS SIMULATORS

58 2.1 | Motivations

59 With respect to available tools for simulating spatially explicit genetic variation across a landscape, SPLATCHE (Currat
60 et al., 2004, 2019) is a user-friendly simulation software that has been supporting the community for two decades,
61 but it is closed source and limited in its configuration capabilities. Considering the wide range of systems that can
62 potentially be analyzed using iDDC, there is not a one-size-fits-all solution: different systems will inevitably require
63 different sets of assumptions/models/simulators. This is exemplified by the many modified versions of the program
64 SPLATCHE used across the literature (e.g., White et al., 2013; Mona et al., 2014). Moreover, because the code is
65 closed source, modifications are restricted to a limited number of people who work with the program and their ability
66 (and availability) to incorporate new implementations.

67 To encourage the open-source creation of new simulators and foster the analysis of new biological systems,
68 QUETZAL-COATL (Becheler et al., 2019; Becheler and Knowles, 2020) was designed as a C++ library of generic
69 components that can be programmed and assembled into versatile simulators. However, its use is by definition re-
70 stricted to C++ programmers (although online tutorials may shorten the beginners learning curve). To widen the range
71 of models available to non-programmers, and incorporate information about the landscape for informing the spatial
72 coalescent, we introduce the open-source QUETZAL-EGGS (<https://github.com/Becheler/QUETZAL-EGGS>).

73 QUETZAL-EGGS contains ready-to-use simulators for implementing different variants of iDDC models. For ex-
74 ample, EGG1 has been developed to simulate fine-grain spatial structure in a system of continental islands formed
75 by progressive submersion of the continental shelf as a response to sea level change after the LGM, but whose pop-
76 ulations remain connected to the mainland by transient trans-oceanic dispersal (e.g., rafting events), whereas EGG2
77 has been developed to illustrate climate-driven pulses in matrix connectivity among relatively isolated populations,
78 such as among montane sky-islands systems (e.g., climate-induced elevational distribution shifts). QUETZAL-EGGS
79 programs take as general inputs a configuration file, a geospatial file describing the landscape of interest and its dy-
80 namics (generally a suitability raster from an ENM step, or multiple rasters for ENMs from different time periods), and
81 a table of sampling locations (latitudinal and longitudinal coordinates). QUETZAL-EGGS complements other spatial
82 simulation resources (e.g., slendr Petr et al., 2022) by offering a compromise between model complexity and compu-
83 tational efficiency. For example, slendr and its backend SLIM (Haller and Messer, 2019) have features to represent
84 spatial interactions between individuals, but the demographic events have to be compiled into a R object, which is
85 expected to be computationally challenging when countless migrations events happen across a complex landscape
86 during a long period of time, compared with SPLATCHE3 and QUETZAL-EGGS simulators that are compiled in C++,
87 and as such, extends the model/parameters space for spatial simulations. We again note that anyone is welcome to
88 contribute to discussions, or can update and grow this list of historical scenarios by adding new models using the
89 Github *Issues* or *Pull Request* systems, or by contacting the authors.

90 2.2 | Memory management

91 One of the significant improvements with our QUEZTAL pipeline regards the computational expense of iDDC mod-
92 eling. For example, SPLATCHE (Currat et al., 2019) keeps the demographic history on RAM, and as such, individual
93 simulations run faster. However, this comes at a cost of constraining the historical duration and landscape resolution

94 (i.e., number of demes) to the system RAM capacity. Because RAM is a more limited resource than disk space, this
95 constrains the number of nodes one can request on computing grids, slowing down the whole workflow and leading
96 to very long run times. In response, researchers try to bypass this problem by re-scaling generation time and/or us-
97 ing coarser landscapes (i.e. to reduce the number of generations and/or number of demes in the spatial simulation,
98 respectively, see He et al., 2013), but this makes other parameters of the model difficult to interpret (Massatti and
99 Knowles, 2016) and prevents the emergence of a fine-grain genetic structure that is often a desirable property for
100 hypothesis testing. To mitigate computational constraints, QUETZAL-EGGS offers a compile-time option that imple-
101 ments sliding windows that keeps only two active layers (i.e., two generations of the spatially explicit demographic
102 history informed by environmental heterogeneity) on RAM at a time, storing unused layers on disk. This allows longer
103 histories at higher spatial resolutions to be modelled.

104 3 | QUETZAL-CRUMBS: PYTHON COMPONENTS SUPPORTING QUETZAL- 105 EGGS

106 3.1 | Motivations

107 A number of iDDC related procedures are not *per se* the responsibility of the simulation program and would require
108 some coding expertise to implement. Below we describe a new python3 library, QUETZAL-CRUMBS, that gathers pro-
109 cedures of general interest for iDDC modeling using QUETZAL-EGGS, improving the accessibility of iDDC modelling
110 to a broad user base.

111 3.2 | Visualization of dynamics landscapes

112 An important part of model choice and calibration is to visually investigate the landscape historical dynamics, whether
113 it is how the candidate model and its parameters affect the demographic history, or how the suitability landscape
114 changes through time. These 2D quantities are represented at each time step by a geospatial regular grid associated
115 to a Coordinate Reference System (a raster). The temporal heterogeneity is represented by stacking these rasters (a
116 multiband raster), where each layer (or band) represents a landscape at a given time period. To visualize how these
117 stacks change through time, the `crumbs.animate` function converts these stacks into GIF or MP4 animations.

118 3.3 | Preparing the landscape and adjusting the spatial grid properties

119 In spatial dynamic models, resolution of the landscape is an issue (see e.g., Bocedi et al., 2012): if the resolution is too
120 low (i.e., large environmentally heterogeneous geographic areas represented as a single deme), biological processes
121 may be misrepresented and biases may result. If the landscape resolution is too high, computational costs may make
122 ABC methodology impossible. Likewise, orientation of the spatial grid is a necessary model parameter, but with
123 multiple orientations possible, this decision is made arbitrarily. To deal with these uncertainties, a common practice is
124 to arbitrarily set a North-up orientation for the spatial grid, and manually guess and adjust the landscape grid resolution
125 to fit computational capacities. However, the impact on inference should be carefully assessed and one way to do so
126 is to include the spatial resolution and grid orientation as parameters to be estimated (e.g., Baird and Santos, 2010;
127 Estoup et al., 2010).

128 The `crumbs.rotate_and_rescale` function allows the rejection of a sample rotation angle/resolution that can
129 not account for the genetic structure of an empirical data set or the simulation walltime is reached. That is, with this

130 QUEZTAL-CRUMBS function, the user can avoid too coarse or too fine of spatial landscape grids and identify the
131 rotation angle that provides the best fit to the observed geographic distribution of genetic variation.

132 3.4 | Beyond the squared spatial grid

133 There are many ways to discretize (tesselate) a landscape. There has been a focus on discrete grids for iDDC modeling
134 partly because SPLATCHE relies on ASCII raster format. However, it is expected that different tessellation models
135 could affect the inference (Baird and Santos, 2010), and consequently, they should be tested. Moreover, considering
136 different tessellations would allow an efficient integration of key data and processes that operate at different scales,
137 such as capturing local micro-refugia without paying the cost of a landscape-wide high resolution (see e.g., Larsson
138 et al., 2021; Randin et al., 2009; Trivedi et al., 2008). Since QUETZAL-COATL embeds abstract libraries like GDAL,
139 the module does not make strong assumptions about tessellation models, requiring only a concept of coordinates,
140 vicinity and distance for sampled individuals/populations. Consequently, different functions to discretize space (like
141 Voronoi tessellations) can be investigated using QUETZAL-CRUMBS; the shapefiles would then be passed on to the
142 QUETZAL-EGGS simulator.

143 Rectangular landscapes can have counter-intuitive orientations that are not very convenient to work with, when
144 compared to disk (circular) landscapes. To facilitate landscape manipulation and analysis, we implement a function
145 `circle_mask` that fits and cuts a circle with maximal radius around the landscape center coordinate when rotating
146 and re-scaling landscapes.

147 3.5 | Representation of temporal heterogeneity at fine scales

148 Despite appreciable progress in accounting for spatial heterogeneity, iDDC studies have focused on a limited number
149 of bands (that is, raster layers) to represent temporal variability (e.g., 1 for static ENM, 2 or 3 for dynamic ENM,
150 see He et al., 2013). This in large part reflects limitations with the available tools for spatially explicit modeling across
151 temporally varying landscapes (Larsson et al., 2021), without some scripting required (e.g., Brown and Knowles, 2012).

152 To ease this step, the `crumbs.interpolate` function takes a n -bands GeoTiff and assigns its first band to genera-
153 tion 0 and its last n band to the simulation maximal generation parameter g (that is, the present). The $n - 2$ remaining
154 bands are then assigned to generations in a regular sequence $[0..g]$, or to a specific sequence provided by the user.
155 Using `dask` (Rocklin, 2015) for parallel computing and larger-than-memory data management, the whole spatial dy-
156 namics is reconstructed by interpolating the missing bands (i.e., bands without independent paleoclimatic data; (see
157 Brown and Knowles, 2012), and this temporal heterogeneity can be animated using `crumbs.animate` and passed to
158 a QUETZAL-EGGS simulator for simulating g generations of a spatial dynamic across the landscape. Note that recon-
159 structing a suitability band for every generation may not scale well to the case of long histories in large landscapes. In
160 these cases, the GDAL Virtual Format (VRT driver, `.vrt`) can be used to build a virtual dataset composed from other
161 GDAL datasets with re-positioning; this allows for very large datasets where most of the bands are actually repeated
162 and reused, rather than physically represented in memory.

163 Rather than interpolating temporal heterogeneity from a few reference paleoclimatic ENMs for iDDC modeling
164 (e.g., Knowles and Massitti 2018), the CHELSA-Trace21k database (Karger et al., 2016) offers high resolution spatio-
165 temporal reconstructions for bioclimatic and elevational data for every century from the present to the LGM (that
166 is, 220 time steps, with a band each 100 years). Using the `crumbs.get_chelsa` function in QUETZAL-CRUMBS, the
167 database variables are downloaded with a procedure that clips and assembles the the 220 layers into a GeoTiff dataset
168 for the spatial extent of the sampled data points (the user can specify a margin to extend the landscape to the desired

size). This automation reduces memory usage and the resultant GeoTiff datasets can be processed by other QUETZAL-CRUMBS modules and by the QUETZAL-EGGS simulators. Note that the long download step can easily be distributed on cluster grids.

Together these advances provide a seamless, flexible iDDC workflow that is also open to extensions. Specifically ease of the iDDC workflow is made possible by (i) databases with major past climatic reconstructions (e.g., Worldclim, Fick and Hijmans, 2017; PaleoClim, Brown et al., 2018; CHELSA, Karger et al., 2016), (ii) ENM software tools (e.g., SDMToolbox Brown, 2014; the R `dismo` package, Hijmans et al., 2017), and (iii) the QUETZAL iDDC modeling framework (Becheler et al., 2019) that generate and/or accept user provided GeoTIFF files.

3.6 | Automated High resolution SDM reconstruction

Using a shapefile of sampling coordinates, the `crumbs.sdm` module fetches CHELSA-Trace21k layers, crops them to the area of interest, and performs a species distribution reconstruction by automatically fetching presence points from the GBIF database (or user input files of occurrences) using 4 machine-learning classifiers (namely, Random Forest, Extra Trees, XGB and LGBM classifiers) to perform model fitting with a k-fold cross validation for computing accuracy scores. The models are then projected to past CHELSA-Trace21k layers and a geotiff is assembled.

We are aware of the numerous challenges that SDM involves and debate regarding the best way to generate SDMs. Here we traded heavily customized approaches for a more general and reproducible workflow. This enables non-programmers to produce a suitability layer for every century during the last 21,000 years, and supply these 220 layers of spatial dynamics to the QUETZAL-EGGS genetic simulation programs. Despite errors with the suitability predictions that may result from this more general and simplified automation of modeling landscape suitability, the inferred suitability predictions can be transformed by an arbitrary function (whose parameters can be estimated by ABC) to improve the generative fit of a model to observed genetic variation. This approach is adopted here because for time periods in the more distant past, assumptions for generating highly precise and accurate projections may not hold (e.g., niche conservatism and or similar community composition such that the species interactions are stable and therefore the relationship between specific environmental variables and a species distributions does not change over time). This contrasts with practices for short-term (the present or decadal) predictions where a highly-precise model may be desirable.

3.7 | Genetic simulation and conversion tools

Because in their current version QUETZAL-EGGS simulate coalescent trees in a Newick format that is stored in a SQLITE database along simulation parameters, QUETZAL-CRUMBS handles access to the simulation SQLITE database, and includes simulation of independent DNA sequences (using Pyvolve; Spielman and Wilke, 2015), data format conversion, and summary statistics computation (using Arlsumstat; Excoffier and Lischer, 2010).

Note that for parameter estimation, QUETZAL-CRUMBS implements procedures already covered by pre-existing libraries (Wegmann et al., 2010; Mertens et al., 2018) to simplify bash scripting and dependency management for genetic simulation under specified priors.

In addition, for spatially explicit simulations, initialization of the simulations has a geographic component. Sometimes the geographic origin might be specified (e.g., based on the putative location of glacial refugia; see Bemmels et al., 2019). However, in other situations the origin is unknown and has to be inferred (He et al., 2017); the `crumbs.sample` function randomly samples candidate origin coordinates among the terrestrial cells of a landscape file in geoTIFF format.

208 3.8 | Sensitivity of inference to sampling of individuals

209 Practical constraints may affect the sampling of individuals across a landscape (e.g., costs of genotyping many individ-
210 uals or difficulties with being able to collect specimens). However, the sampling scheme itself may impact inferences
211 made from genetic data (Mason et al., 2020). For example, limited sampling of geographically widespread taxa may
212 generate genetic patterns that deviate from coalescent expectations for a single population, and as a consequence,
213 the data might fit better a "multispecies" coalescent, MSC, model (i.e., more than one population lineage). In such
214 cases, the sampling (rather than limited gene flow) would drive support for multiple population lineages, which in
215 turn, is commonly interpreted as support for multiple cryptic species in the parlance of species delimitation (Barley
216 et al., 2018; Sukumaran and Knowles, 2017). Yet, tests for such biased inferences arising from the sampling design
217 are not common.

218 The program DECRYPT, which uses QUEZTAL-EGGs to simulate a spatial coalescent informed by the environment
219 (i.e., habitat suitability using QUEZTAL-CRUMBS), can be used to test for sensitivities due to sampling. Specifically,
220 simulated data sets from the posterior of a full iDDC model (i.e., pseudo-observed data sets, PODs) are used to
221 assess the robustness of the MSC to possible violation of its assumptions (e.g., restrictions in gene flow arising from
222 environmental heterogeneity). That is, for a particular geographic sample design (the actual geographic coordinates
223 of empirical samples) the inferred number of lineages can be estimated under the MSC. This provides a test of the
224 robustness of the MSC to violations of the models assumptions, such as genetic structure within a species as an
225 artifact of the sampling scheme or due to reduced gene flow because of landscape features.

226 4 | QUETZAL-NEST

227 For a non-programmer and newcomer to iDDC modeling, one the first barrier encountered is the diversity and dis-
228 persion of tools and methods: identifying, installing, configuring, calibrating and running the required tools is far from
229 trivial, even for simple tests on a local computer, and not to mention runs conducted on a cluster for scalable, repro-
230 ducible science. A streamlined software solution that alleviates at least some of the complexity of analyses based on
231 a spatial coalescent model is key to broadening the scope of potential users.

232 Ideally, we will see the emergence of a framework for reproducible iDDC where the practitioner would only
233 have to (i) connect to an HTC grid, (ii) download content from a standardized Github repository of gathered tools
234 and methods for analysis, (iii) upload their own input files, (iv) select and run routine analysis workflows, and (v)
235 retrieve and interpret outputs. Recent advances make some of this path a bit easier. First, the recent developments
236 of ABC-Random Forest (Raynal et al., 2019) now allow scientists to perform ABC inference, bypassing complex and
237 time-consuming aspects of the inference, which enables the design of more standard ABC workflows. Second, the
238 emergence of containers (e.g., Docker and Singularity Kurtzer et al., 2017) and distributed High Throughput Computing
239 (e.g., the Open Science Grid, Pordes et al., 2007) now allow packages to be shared and run in reproducible analytical
240 environments.

241 As a first step in this direction, we developed the QUETZAL-NEST Docker container that comes with about 65 pre-
242 installed dependencies. The container is published on DockerHub and available for local use (e.g., development, tests,
243 tutorials) with `docker pull arnaudbecheler/quetzal-nest`. To allow researchers to perform full iDDC inferences
244 with ABC, QUETZAL-NEST has also been submitted to the Open Science Grid CVMFS image repository where it is
245 available for distributed High Throughput Computing. An example repository of a full data analysis workflow built for
246 OSG can be found at https://github.com/Becheler/quetzal_on_OSG.

5 | FUTURE PROSPECTS

Going forward, landscape types, demographic and historical details, and geographic settings will expand beyond the current resources of QUEZTAL, and will be made available as additional QUEZTAL-EGGS simulators beyond the current list. Such additions are eased by the clear structure and intent of the C++ files that define each existing QUEZTAL-EGGS: `EGG_options.h` defines the simulator options, `EGG_context.h` defines the forward/backward model, `EGG_database` is responsible for storing simulated parameter values and data, and the `main.cpp` contains the main function. All files are relatively short and the code can be modified with minimal C++ knowledge.

Currently, Quetzal simulates independent loci. Although this assumption is simple, it still matches a large number of existing geospatial genetic datasets. However, it also ignores the rich information embedded in recombination patterns of more complex datasets. Given its open source code and abstract interfaces, Quetzal could be interfaced with TSKIT (Kelleher et al., 2018) and/or SLIM for computationally efficient generated spatial history of whole genomes. More specifically, we began to implement a C++ version of the Hudson algorithm for enabling the simulation of a structured coalescent with recombination.

With this flexibility in mind, we have developed QUEZTAL so it can continue to evolve to fit future demands of spatially explicit genetic studies in an open environment that is available to all researchers.

acknowledgements

We thank two anonymous reviewers for useful comments of an earlier version of the paper. This work was supported by NSF DEB 1655607 (to LLK).

conflict of interest

None declared.

references

- Baird, S. J. and Santos, F. (2010) Monte carlo integration over stepping stone models for spatial genetic inference using approximate bayesian computation. *Molecular ecology resources*, **10**, 873–885.
- Barley, A. J., Brown, J. M. and Thomson, R. C. (2018) Impact of model violations on the inference of species boundaries under the multispecies coalescent. *Systematic Biology*, **67**, 269–284.
- Beaumont, M. A., Zhang, W. and Balding, D. J. (2002) Approximate bayesian computation in population genetics. *Genetics*, **162**, 2025–2035.
- Becheler, A., Coron, C. and Dupas, S. (2019) The quetzal coalescence template library: A c++ programmers resource for integrating distributional, demographic and coalescent models. *Molecular ecology resources*, **19**, 788–793.
- Becheler, A. and Knowles, L. L. (2020) Occupancy spectrum distribution: application for coalescence simulation with generic mergers. *Bioinformatics*, **36**, 3279–3280.
- Bemmels, J. B., Knowles, L. L. and Dick, C. W. (2019) Genomic evidence of survival near ice sheet margins for some, but not all, north american trees. *Proceedings of the National Academy of Sciences*, **116**, 8431–8436.
- Bocedi, G., Pe'er, G., Heikkinen, R. K., Matsinos, Y. and Travis, J. M. (2012) Projecting species' range expansion dynamics: sources of systematic biases when scaling up patterns and processes. *Methods in Ecology and Evolution*, **3**, 1008–1018.

- 282 Brown, J. L. (2014) Sdm toolbox: a python-based gis toolkit for landscape genetic, biogeographic and species distribution
283 model analyses. *Methods in Ecology and Evolution*, **5**, 694–700.
- 284 Brown, J. L., Hill, D. J., Dolan, A. M., Carnaval, A. C. and Haywood, A. M. (2018) Paleoclim, high spatial resolution paleoclimate
285 surfaces for global land areas. *Scientific data*, **5**, 1–9.
- 286 Brown, J. L. and Knowles, L. L. (2012) Spatially explicit models of dynamic histories: examination of the genetic consequences
287 of pleistocene glaciation and recent climate change on the american pika. *Molecular Ecology*, **21**, 3757–3775.
- 288 Currat, M., Arenas, M., Quilodràn, C. S., Excoffier, L. and Ray, N. (2019) Splat3: simulation of serial genetic data under
289 spatially explicit evolutionary scenarios including long-distance dispersal. *Bioinformatics*, **35**, 4480–4483.
- 290 Currat, M., Ray, N. and Excoffier, L. (2004) Splat3: a program to simulate genetic diversity taking into account environmental
291 heterogeneity. *Molecular Ecology Notes*, **4**, 139–142.
- 292 Estoup, A., Baird, S. J., Ray, N., Currat, M., CORNUET, J.-M., Santos, F., Beaumont, M. A. and Excoffier, L. (2010) Combining
293 genetic, historical and geographical data to reconstruct the dynamics of bioinvasions: application to the cane toad *bufo*
294 *marinus*. *Molecular ecology resources*, **10**, 886–901.
- 295 Excoffier, L. and Lischer, H. E. (2010) Arlequin suite ver 3.5: a new series of programs to perform population genetics analyses
296 under linux and windows. *Molecular ecology resources*, **10**, 564–567.
- 297 Fick, S. E. and Hijmans, R. J. (2017) Worldclim 2: new 1-km spatial resolution climate surfaces for global land areas. *International*
298 *journal of climatology*, **37**, 4302–4315.
- 299 Haller, B. C. and Messer, P. W. (2019) Slim 3: forward genetic simulations beyond the wright–fisher model. *Molecular biology*
300 *and evolution*, **36**, 632–637.
- 301 He, Q., Edwards, D. L. and Knowles, L. L. (2013) Integrative testing of how environments from the past to the present shape
302 genetic structure across landscapes. *Evolution*, **67**, 3386–3402.
- 303 He, Q., Prado, J. R. and Knowles, L. L. (2017) Inferring the geographic origin of a range expansion: Latitudinal and longitudinal
304 coordinates inferred from genomic data in an abc framework with the program x-origin. *Molecular Ecology*, **26**, 6908–6920.
- 305 Hijmans, R. J., Phillips, S., Leathwick, J., Elith, J. and Hijmans, M. R. J. (2017) Package 'dismo'. *Circles*, **9**, 1–68.
- 306 Karger, D. N., Conrad, O., Böhrner, J., Kawohl, T., Kreft, H., Soria-Auza, R. W., Zimmermann, N. E., Linder, H. P. and Kessler, M.
307 (2016) Chelsa climatologies at high resolution for the earth's land surface areas (version 1.0).
- 308 Kelleher, J., Thornton, K. R., Ashander, J. and Ralph, P. L. (2018) Efficient pedigree recording for fast population genetics
309 simulation. *PLoS computational biology*, **14**, e1006581.
- 310 Knowles, L. L. and Alvarado-Serrano, D. (2010) Exploring the population genetic consequences of the colonization process
311 with spatio-temporally explicit models: insights from coupled ecological, demographic and genetic models in montane
312 grasshoppers. *Molecular Ecology*, **19**, 3727–3745.
- 313 Kurtzer, G. M., Sochat, V. and Bauer, M. W. (2017) Singularity: Scientific containers for mobility of compute. *PLoS one*, **12**,
314 e0177459.
- 315 Larsson, D. J., Pan, D. and Schneeweiss, G. M. (2021) Addressing alpine plant phylogeography using integrative distributional,
316 demographic and coalescent modeling. *Alpine Botany*, 1–15.
- 317 Mason, N. A., Fletcher, N. K., Gill, B. A., Funk, W. C. and Zamudio, K. R. (2020) Coalescent-based species delimitation is
318 sensitive to geographic sampling and isolation by distance. *Systematics and Biodiversity*, **18**, 269–280.
- 319 Massatti, R. and Knowles, L. L. (2016) Contrasting support for alternative models of genomic variation based on microhabitat
320 preference: Species-specific effects of climate change in alpine sedges. *Molecular Ecology*, **25**, 3974–3986.

- 321 Mertens, U. K., Voss, A. and Radev, S. (2018) Abrox—a user-friendly python module for approximate bayesian computation
322 with a focus on model comparison. *PLoS one*, **13**, e0193981.
- 323 Mona, S., Ray, N., Arenas, M. and Excoffier, L. (2014) Genetic consequences of habitat fragmentation during a range expansion.
324 *Heredity*, **112**, 291–299.
- 325 Neuenschwander, S., Lurgiader, C. R., Ray, N., Currat, M., Vonlanthen, P. and Excoffier, L. (2008) Colonization history of the
326 swiss rhine basin by the bullhead (*cottus gobio*): inference under a bayesian spatially explicit framework. *Molecular Ecology*,
327 **17**, 757–772.
- 328 Ortego, J. and Knowles, L. L. (2020) Incorporating interspecific interactions into phylogeographic models: A case study with
329 californian oaks. *Molecular Ecology*, **29**, 4510–4524.
- 330 Pan, D., Hülber, K., Willner, W. and Schneeweiss, G. M. (2020) An explicit test of pleistocene survival in peripheral versus
331 nunatak refugia in two high mountain plant species. *Molecular ecology*, **29**, 172–183.
- 332 Petr, M., Haller, B. C., Ralph, P. L. and Racimo, F. (2022) slendr: a framework for spatio-temporal pop-
333 ulation genomic simulations on geographic landscapes. *bioRxiv*. URL: [https://www.biorxiv.org/
334 content/early/2022/03/21/2022.03.20.485041](https://www.biorxiv.org/content/early/2022/03/21/2022.03.20.485041). Publisher: Cold Spring Harbor Laboratory _eprint:
335 <https://www.biorxiv.org/content/early/2022/03/21/2022.03.20.485041.full.pdf>.
- 336 Pordes, R., Petravick, D., Kramer, B., Olson, D., Livny, M., Roy, A., Avery, P., Blackburn, K., Wenaus, T., Würthwein, F. et al.
337 (2007) The open science grid. In *Journal of Physics: Conference Series*, vol. 78, 012057. IOP Publishing.
- 338 Randin, C. F., Engler, R., Normand, S., Zappa, M., Zimmermann, N. E., Pearman, P. B., Vittoz, P., Thuiller, W. and Guisan, A.
339 (2009) Climate change and plant distribution: local models predict high-elevation persistence. *Global Change Biology*, **15**,
340 1557–1569.
- 341 Raynal, L., Marin, J.-M., Pudlo, P., Ribatet, M., Robert, C. P. and Estoup, A. (2019) Abc random forests for bayesian parameter
342 inference. *Bioinformatics*, **35**, 1720–1728.
- 343 Rocklin, M. (2015) Dask: Parallel computation with blocked algorithms and task scheduling. In *Proceedings of the 14th python
344 in science conference*, vol. 130, 136. Citeseer.
- 345 Spielman, S. J. and Wilke, C. O. (2015) Pyvolve: a flexible python module for simulating sequences along phylogenies. *PLoS
346 one*, **10**, e0139047.
- 347 Sukumaran, J. and Knowles, L. L. (2017) Multispecies coalescent delimits structure, not species. *Proceedings of the National
348 Academy of Sciences*, **114**, 1607–1612.
- 349 Trivedi, M. R., Berry, P. M., Morecroft, M. D. and Dawson, T. P. (2008) Spatial scale affects bioclimate model projections of
350 climate change impacts on mountain plants. *Global change biology*, **14**, 1089–1103.
- 351 Wegmann, D., Leuenberger, C., Neuenschwander, S. and Excoffier, L. (2010) Abctoolbox: a versatile toolkit for approximate
352 bayesian computations. *BMC bioinformatics*, **11**, 1–7.
- 353 White, T. A., Perkins, S. E., Heckel, G. and Searle, J. B. (2013) Adaptive evolution during an ongoing range expansion: the
354 invasive bank vole (*Myodes glareolus*) in Ireland. *Molecular Ecology*, **22**, 2971–2985.

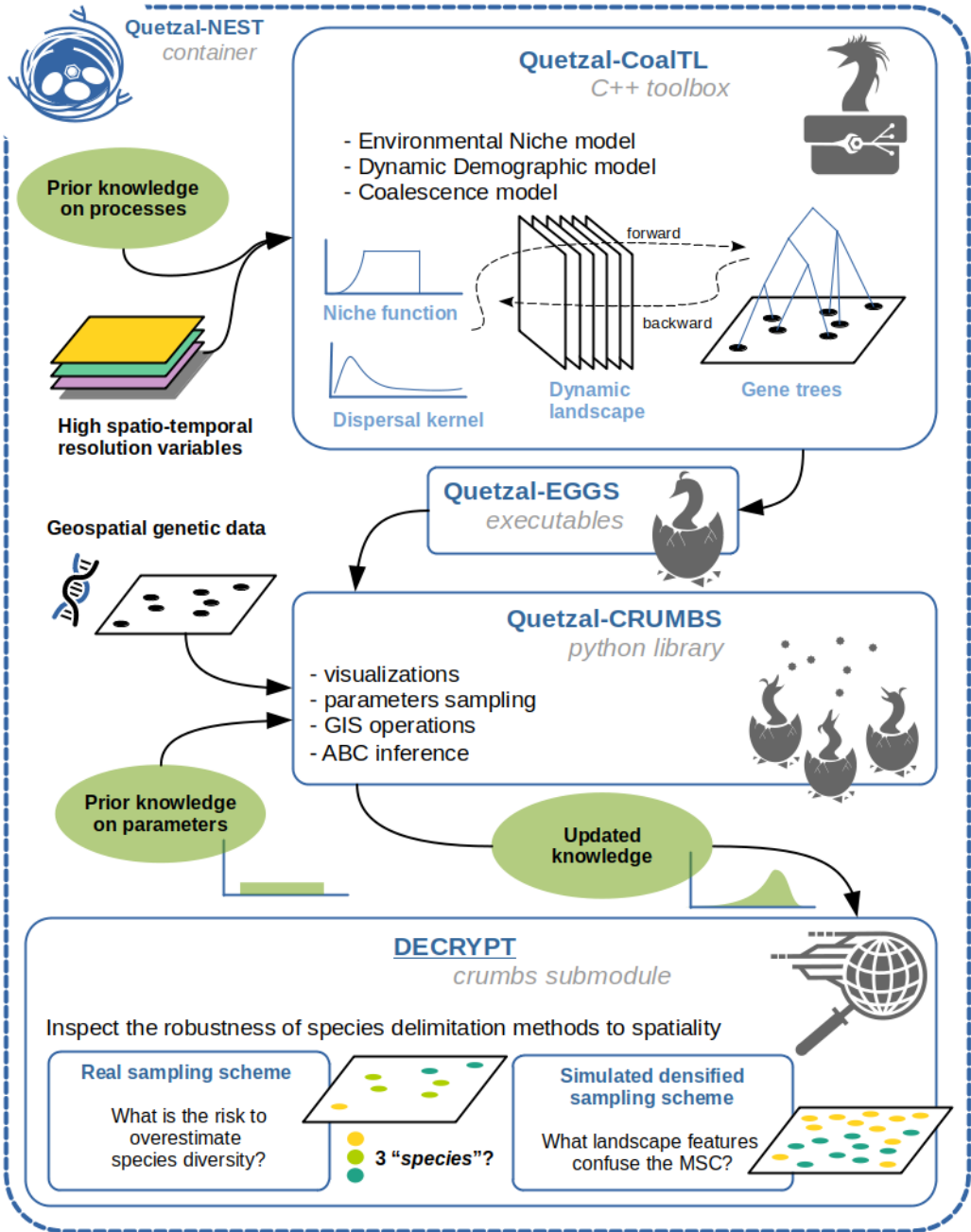


FIGURE 1 Main components, concepts and uses of the QUETZAL framework for open source iDDC modeling. QUETZAL-NEST is a Docker container that packages all the tools and dependencies; it can be run locally with Docker or on dHTC clusters with Singularity. QUETZAL-COATL (Becheler et al., 2019; Becheler and Knowles, 2020) is a C++ library of reusable components and QUETZAL-EGGS are C++ iDDC simulators built with these components. QUETZAL-CRUMBS is a complementary set of Python tools for hypothesis testing using ABC and common landscape-ABC problems, including automatic adjustment of the spatial resolution and orientation of the landscape. DECRYPT is a submodule of CRUMBS for automated, biology-informed robustness analysis of the multispecies coalescent model.