# Probabilistic graph models for landscape genetics

**Brook G. Milligan**[1]

[1]**Department of Biology, New Mexico State University, Las Cruces, New Mexico 88003 USA**

Corresponding author:
Brook G. Milligan[1]

Email address: brook@nmsu.edu

## ABSTRACT

Landscape genetics combines population genetics, landscape ecology, and spatial analysis to identify landscape and genetic factors that influence genetic and genomic variation. Progress in the field depends on a strong conceptual foundation and the means of identifying mechanistic connnections between environmental factors, landscape features, and genetic or genomic variation. Many existing approaches and much of the software commonly in use was developed for population genetics or statistics and is not entirely appropriate for landscape genetics. Probabilistic graph models provide a statistically rigorous and flexible means of constructing models directly applicable to landscape genetics. Probabilistic graph models also allow construction of mechanistic models, which are crucial elements in testing hypotheses. Sophisticated software exists for the analysis of graph models; however, much of it does not handle the types of data used for landscape genetics, model structures involving autoregressive spatial interaction between variables, or the scale of landscape genetics problems. Thus, an important priority for the field is to develop suitably flexible software tools for graph models that overcome these problems and allow landscape geneticists to explore meaningfully mechanistic and flexible models. We are developing such a library and applying it to examples in landscape genetics.

Keywords:    landscape genetics, population genetics, graph models, Bayesian inference, open source software, software development

Landscape genetics combines population genetics, landscape ecology, and spatial analysis to identify the mechanisms by which landscape and environmental factors influence genetic and genomic variation. From the outset, the field has focused on the twin ecological and evolutionary processes of gene flow and adaptation (Holderegger et al., 2006; Manel et al., 2003, 2010). Involving as it does quantification of both genetics and landscapes, landscape genetics is inherently interdisciplinary (Balkenhol et al., 2009; Holderegger and Wagner, 2008). While the emphasis is often on the genetics, explicit consideration of the importance of GIS and allied geospatial disciplines is crucial as they can contribute to landscape genetics in many ways (Cushman et al., 2016; Storfer et al., 2007). For example, experimental design in landscape genetics must be informed by such factors as the spatial extent and grain of available data, and the configuration of landscape features. Landscape and environmental data are inherently spatial, and must be acquired, organized, and analyzed in the course of a landscape genetics study. Thus, geoscientists and geocomputation will play an increasingly important role in landscape genetics.

Progress in landscape genetics is so far limited by available analytical methods (Balkenhol et al., 2009, 2016a; Guillot et al., 2009). In part this derives from the fact that many of the available analytical tools and much of the usable software were originally developed for population genetics or even broader statistical applications. They often include assumptions and are applicable to data that are not completely appropriate for landscape genetics studies. Because of this gap, there is no consensus in the literature regarding how to approach landscape genetics analysis (Balkenhol et al., 2016a). Indeed, the *ad hoc* assortment of methods currently in use lacks a unifying theory; consequently, more focus must be given to a mechanistic understanding of the influence of landscapes and environments on genetic and genomic variation (Balkenhol et al., 2016b). Development of a more comprehensive theory will come in part from an improved foundation of open source computational tools allowing explicit and flexible mechanistic modeling.

This brief review focuses on three themes. First, it identifies the types of models most likely to advance a comprehensive theory of landscape genetics, improve mechanistic understanding, and provide better predictions serving, for example, conservation policy and management. Second, it considers a set of open source software that could be used for general models in landscape genetics but that all have significant limitations. Finally, it also suggests how these limitations can be overcome with new models and computational tools.

# 1 LANDSCAPE GENETICS AND BAYESIAN INFERENCE

The prevailing challenge in landscape genetics is identifying the mechanisms by which landscape and environmental factors influence genetic and genomic variation. More precisely, the central question is: given data on intraspecific genetic variation across landscapes (or waterscapes; Manel and Holderegger (2013); Selkoe et al. (2016)), what inferences are possible regarding the functional mechanisms and factors causing that variation? Framing the question in this way emphasizes the inherent connection between the science of landscape genetics and the nature of Bayesian inference.

The natural connection between landscape genetics and Bayesian inference has led to the development of a variety of widely used Bayesian analysis methods. A first set of these includes STRUCTURE, which identifies putative populations and assigns individuals to them (Pritchard et al., 2000). Although originally designed for population not landscape genetics, it remains the most widely used. A second set of Bayesian models applied to landscape genetics includes GENELAND, which seeks to identify population clusters by modeling allele frequency distributions in a spatially explicit way (Chen et al., 2007; Guillot et al., 2005a,b). More recently, Bayesian models that explicitly relate environmental gradients to spatially explicit allele frequency distributions have been developed (Coop et al., 2010; Frichot et al., 2013).

One element is common to all of the available software: each program implements a narrow range of possible models and provides very limited opportunity for expanding its scope. For example, as discussed below, both STRUCTURE and GENELAND are essentially variants of the same model, yet nothing of their implementation is shared so new variants cannot be created by exploiting their commonality. Further, the published descriptions do not reveal the inherent similarity between STRUCTURE and GENELAND, so conceptual connections are not evident. Consequently, landscape geneticists do not recognize a continuum of possible models. Even worse, they cannot exploit the continuum by incrementally modifying existing models and competing alternatives against available

data. This is a serious limitation for a scientific field that repeatedly asserts that more mechanistic and predictive models and a stronger theoretical foundation are essential (Andrew et al., 2013; Balkenhol et al., 2016b; Guillot et al., 2009; Manel and Holderegger, 2013).

## 2 PROBABILISTIC GRAPH MODELS

Mathematical graphs are widely used to represent models, including some in landscape genetics. Graphs are composed of a set of vertices and a set of edges, each of which connects a pair of vertices. Edges may be directed or undirected, and paths are sequences of edges connecting one vertex with another, possibly with intervening vertices. A cyclic graph has at least one path starting and ending at the same vertex; an acyclic graph lacks any such paths.

One application of graphs to landscape genetics derives from the population graph concept (Dyer and Nason, 2004). Here the graph is composed of vertices representing population distributions in a multilocus genetic space, and edges representing interdependencies between populations due, for example, to gene flow (Excoffier et al., 1992). The primary application to landscape genetics has been identification of conditional independence between populations to remove edges followed by analysis of graph structure metrics such as centrality or connnectness (Dyer, 2007; Murphy et al., 2016).

Graph models can be much richer, however, and both STRUCTURE and GENELAND are examples used in landscape genetics. Generally, (probabilistic) graph models are composed of vertices representing any kind of random variable and edges representing dependencies between them (Bishop, 2006; Koller and Friedman, 2009). They are widely used, for example, in latent factor analysis (Steyvers and Griffiths, 2007), a field that now finds application broadly in machine learning, artificial intelligence, and document and image processing, as well as landscape genetics (Blei et al., 2003; Blei, 2012; Frichot et al., 2013; Jia et al., 2011; Pritchard et al., 2000). The population graph concept of Dyer and Nason (2004) is clearly a special case where each vertex represents the same quantity, a population-specific distribution, but the landscape genetics analysis involving edge removal and graph metrics (Murphy et al., 2016) is unrelated to the use of graphs as formal probabilistic models (Bishop, 2006; Koller and Friedman, 2009). The value of the latter for landscape genetics, both conceptually and for software development, is the focus here.

Although not described as such, a probabilistic graph model represents the mathematics underlying STRUCTURE (Pritchard et al., 2000). In this case, the random variables represent population-specific distributions of alleles, the probabilistic assignment of alleles to populations, and prior distributions that by default are uninformative (Figure 1). The STRUCTURE software supports slight variations in the model depicted; for example, assignment of all alleles may be individual-specific not allele-specific as shown, and priors may be informative in various ways. These variations, however, are extremely limited and do not cover the continuum of related models that is possible.

One related model, however, is alluded to in Pritchard et al. (2000) and described in detail in Falush et al. (2003); but again, the graph model itself is not presented explicitly. The main difference is that in this model the population-specific allele distributions are not independent; instead, they are correlated via a shared ancestral population (Figure 2).

A further related model, implemented in GENELAND, is described in Guillot et al. (2005a), again without depicting the graph model (Figure 3). This model explicitly adds spatial information

to the model; unlike the other two, both the identity of alleles and their spatial location are observed. This supports estimating additional random variables such as the inferred location of individuals and spatially-explicit allele distributions.

A comparison of Figures 1–3 makes clear that these are all closely related models, a fact that is generally not made evident by the papers describing them. Furthermore, in many ways the graph models are more useful than the papers, because they make the conceptual linkages clear and enable direct comparisons among them. They also make gaps in the existing models evident; for example, none of these include gene flow explicitly despite its clear importance as a mechanism in landscape genetics (Holderegger and Wagner, 2008; Manel and Holderegger, 2013; Storfer et al., 2007; van Strien et al., 2014). Finally, probabilistic graph models invite the construction of variations by adding new random variables or changing dependencies among them, because the biological structure of the models is easy to reason about when presented in the form of a graph. Probabilistic graph models, therefore, provide an ideal foundation for mechanistic modeling in landscape genetics that can lead to an improved theoretical understanding.

## 3 A MECHANISTIC MODELING FRAMEWORK FOR LANDSCAPE GE-NETICS

Traditional approaches to landscape genetics descriptively model either genetic characteristics associated with each sampled site or individual, or derived genetic measures associated with pairs of sampled sites or individuals (Joost et al., 2007). Almost all approaches model these response variables using *ad hoc* distributions taken from more generic statistical literature; for example, virtually the entire textbook on landscape genetics (Balkenhol et al., 2016a) follows this pattern. In contrast, a mechanistic approach would construct a model of the individual observations, e.g., individual multilocus genotypes (or genomes), as a function of assumed demographic, ecological, and population genetic mechanisms.

As described earlier and illustrated in Figures 1–3, STRUCTURE and GENELAND are examples of exactly this approach; the observed alleles are modeled directly in terms of unobserved but inferable populations and assignments (Guillot et al., 2005a; Pritchard et al., 2000). Viewed in this context, differences between individual- and population-based approaches to landscape genetics are not fundamental; rather they reduce to simple differences between the structure of the graphical model in use. Individual-based models have graphs that relate observations on individuals to individual-specific random variables; examples of the latter are the assignment of an individual's alleles to populations ($Z$ in Figures 1 and 3) and the inferred true location of each individual ($s$ in Figure 3). Population-based models have graphs that relate observations on individuals to population-specific random variables; examples of these are the population-specific allele frequencies ($P$ in Figures 1 and 3). By including elements of each, Figures 1 and 3 already blur the boundary between individual- and population-specific models.

Given the power of probabilitistic graph models to represent a broad spectrum of intermediate cases just as well, a better framework is the set of mechanisms included. From this perspective, it is evident that Figure 3 includes spatially-explicit mechanisms whereas Figure 1 does not. It is also evident that neither one includes an explicit mechanism for gene flow. The power of probabilitistic graph models lies in their ability to cover the entire spectrum of models relevant to landscape

genetics and to encourage more transparent reasoning about alternative models. Using them to advance landscape genetics is limited only by our ability to compare alternative models, but that in turn is severely constrained by the software available to manipulate and analyze them.

# 4 OPEN-SOURCE PROBABILISTIC GRAPH MODELS

As just illustrated, the primary advantages of probabilistic graph models are that complex and realisticly mechanistic models can be constructed, and that their model structure can be manipulated easily to explore alternatives. Thus, there is great scope for constructing general theories based upon manipulating probabilitistic graph models to reflect interesting biological models within landscape genetics. However, software tools must exist that enable manipulation and analysis of the graphs, and the types of graphs available must match those required by landscape genetics. For many applications two types of graphs are enough: Bayesian networks represented by directed acyclic graphs (DAGs) and Markov random fields represented by undirected graphs. Landscape genetics models, however, often require more general types of graphs to accommodate, for example, spatially autoregressive relationships among random variables. Additionally, landscape genetics models often require distributions appropriate to a broad range of commonly encountered data types, including alleles, genotypes, spatially explicit environmental data. Such a range of discrete and continuous, unidimensional and multidimensional data types requires a rich array of probability distributions.

While the set of probabilistic graph models that has been applied to landscape genetics do not harness their full flexibility, there exist modeling software that does better. The most widely used is based upon the BUGS language for describing graph models, and includes WinBUGS, OpenBugs (Lunn et al., 2009) and JAGS (Plummer, 2015). The BUGS language allows textual description of general graph models that include a broad range of distributions. The textual description is translated into executable code, a process that introduces some of the limitations common to this type of modeling software. First, the flexibility of possible applications is limited by the features of the BUGS language. A limited range of data types, generally scalars and vectors or matrices constructed from them, is available, only data structures describable in the language may be used, and algorithms are limited to those already programmed. Second, the scale of models is also limited by the execution environment provided by the implementation. Despite the inherent flexibility of graph models in general, both of these limitations are barriers to convenient development of landscape genetics models that leverage the flexibility of graph models. While genetic data can be recoded in the form of only integers or real numbers, it is tedious and error-prone to do so; thus, the limited data types available create needless barriers. A landscape genetics model might include thousands or millions of random variables within it; consider, for example, a model of population allele freqencies and environmental factors across a landscape grid of $1000 \times 1000$ pixels. This puts severe stress on models that cannot harness the full power of multithreading, distributed multiprocessing, and careful memory management. Being limited by the BUGS language, these programs provide restricted capacity for modelers to address these issues.

Another general graph modeling system is Stan (Carpenter et al., 2015; Gelman et al., 2015). Although more flexible in some ways than BUGS, Stan suffers from some of the same limitations that reduce its applicability to landscape genetics. It has the same limited data types and the execution environment is likewise limited by the Stan language. As a result, neither BUGS nor Stan are ideally

| Name | Graph types | Primitive variables | Preprocessing | Implementation language | Reference |
|------|-------------|---------------------|---------------|------------------------|-----------|
| Darwin | FGs | scalars | compiled | C++ | Gould (2015) |
| HYDRA | DAGs, MRFs, FGs, HMMs | Java classes | compiled | Java | Warmes (2013) |
| Infer.NET | FGs | C# classes | compiled | C# | Minka et al. (2014) |
| JAGS | DAGs | scalars | interpreted | C++ | Plummer (2016) |
| JavaBayes | DAGs | scalars | interpreted | Java | Cozman (2001) |
| libDAI | FGs | discrete | compiled | C++ | Mooiji (2015) |
| Mocapy++ | DAGs, HMMs | C++ classes | compiled | C++ | Antonov et al. (2015) |
| Nimble | DAGs | scalar | interpreted | C++ | de Valpine et al. (2016) |
| OpenBUGS | DAGs | scalar | interpreted | Component Pascal | Thomas (2009) |
| OpenGM | DAGs, MRFs, FGs | discrete | compiled | C++ | OpenGM (2015) |
| PNL | DAGs, MRFs | C++ classes | compiled | C++ | Sysoyev et al. (2013) |
| RISO | DAGs | Java classes | compiled | Java | Dodier (2012) |
| Stan | | scalars | interpreted | C++ | Stan Development Team (2016) |
| Vibes | DAGs | scalar | compiled | Java | Winn (2004) |

**Table 1.** A selection of open source software tools for analyzing probabilistic graph models. Type of graphs include directed acyclic graphs (DAGs), Markov random fields (MRFs), factor graphs (FGs), hidden Markov models (HMMs), and Gaussian Markov models (GMMs).

suited for landscape genetics applications.

In addition to these two major classes of graph modeling software, a broad range of more specialized software systems is also available; many of these are summarized by Murphy (2014). Some are open source and may have potential for landscape genetics applications (Table 1). These tools have many of the same limitations as BUGS, JAGS, and Stan. They often handle fewer graph types than needed for landscape genetics, the data types are not well suited to landscape genetics, or their execution environments are restrictive. In addition, they are much more specialized, difficult to program, and likely well beyond the reach of typical landscape geneticists. These characteristics mean that landscape geneticists face a fundamental challenge hindering development of a strong conceptual foundation for the field based upon the expressive power, flexibility, and statistical rigor of probabilistic graph models.

# 5 DESIGNING A PROBABILISTIC GRAPH MODEL FOR LANDSCAPE GENETICS

What then is the ideal design of a software system intended to harness the power, flexibility, and rigor of probabilistic graph models applied to landscape genetics? First and foremost, it must support a full range of relevant graph types, which in particular means not being limited to directed acyclic graphs. Second, it must support a full range of useful data types that landscape geneticists work with; in addition to simple scalars, vectors, and matrices, these include named alleles and genotypes, loci and chromosomes, geographic locations, and spatial data of various sorts. Ideally, user-defined or third-party data types should be easy to accommodate. Third, the algorithms available should be extensible to allow improved efficiency as needed. Fourth, the execution environment should not be limited to that encapsulated within a single, predefined program. This is especially important for landscape genetics models that may well encompass thousands or millions of random variables. Finally, the power and flexibility of graph models must be abstracted enough that a full spectrum

²⁰⁹ of landscape geneticists can create simple models easily, test alternative and biologically relevant
²¹⁰ models quickly, and improve upon the models and algorithms as needed.

²¹¹     It is little surprise that existing software tools are unable to meet these stringent demands; they
²¹² are largely conflicting and impossible to resolve without advanced software design. The most likely
²¹³ path forward (Lunn et al., 2009) leverages the power of C++ to present high-level abstractions
²¹⁴ based upon embedded domain specific languages (de Guzman and Kaiser, 2016; Niebler, 2016)
²¹⁵ assembled with expression templates (Niebler, 2016; Veldhuizen, 1995) from highly reusable generic
²¹⁶ components (Stepanov and Rose, 2014). Although beyond the scope of this paper, we are following
²¹⁷ these design principles to implement a software library intended to provide the expressive power and
²¹⁸ computational performance demanded for advancing a coherent conceptual foundation for landscape
²¹⁹ genetics. The outcome is a highly compact way of encoding probabilistic graph models of relevance
²²⁰ to landscape genetics (Figure 4). Given the expressive power of the language, all of this should
²²¹ be readily accessible to biologists without deep knowledge of C++ programming. Importantly,
²²² models can be described in a formal way that removes the ambiguity inherent in natural language
²²³ descriptions. Finally, because models are encoded directly in C++, not interpreted, they can be
²²⁴ reused as portions of larger programs for enhanced capability; this is fundamentally impossible for
²²⁵ interpreted modeling frameworks such as OpenBUGS or JAGS. The generality of this approach
²²⁶ removes the limitations inherent to the available software and characteristic of current approaches to
²²⁷ landscape genetics data analysis, and makes it easy to encode, and therefore explore, the complete
²²⁸ space of relevant models.

## 6 CONCLUSION

²³⁰ Landscape genetics suffers greatly from the absence of an analytical foundation that encourages
²³¹ development of a mechanistic understanding of the impact of environmental and landscape factors
²³² on genetic and genomic variation (Balkenhol et al., 2016a). This stems in part from the adoption of
²³³ software tools and methods originally developed for other purposes. There exist well-established
²³⁴ concepts and statistical approaches associated with probabilitistic graph models that are ideally
²³⁵ suited as the needed foundation for landscape genetics. Unfortunately, the associated software tools
²³⁶ cannot be borrowed directly, because they are limited in ways that do not accommodate the needs of
²³⁷ landscape geneticists. One priority that would directly advance the field and resolve these problems
²³⁸ is the development of probabilistic graph model tools that do apply generally to landscape genetics.
²³⁹ Despite the inherent difficulty of this task, we have developed a suitable library and are beginning to
²⁴⁰ apply it to landscape genetics.

## 7 ACKNOWLEDGEMENTS

# REFERENCES

Andrew, R. L., Bernatchez, L., Bonin, A. L., Buerkle, C. A., Carstens, B. C., Emerson, B. C., Garant, D., Giraud, T., Kane, N. C., Rogers, S. M., Slate, J., Smith, H., Sork, V. L., Stone, G. N., Vines, T. H., Waits, L., Widmer, A., and Rieseberg, L. H. (2013). A road map for molecular ecology. *Molecular Ecology*, 22:2605–2626.

Antonov, L., Paluszewski, M., and Hamelrcyk, T. (2015). Mocapy++. `https://sourceforge.net/projects/mocapy/`.

Balkenhol, N., Cushman, S. A., Storfer, A. T., and Waits, L. P., editors (2016a). *Landscape genetics: concepts, methods, applications*. Wiley Blackwell, Hoboken, New Jersey.

Balkenhol, N., Cushman, S. A., Waits, L. P., and Storfer, A. (2016b). Current status, future opportunities, and remaining challenges in landscape genetics. In Balkenhol, N., Cushman, S. A., Storfer, A. T., and Waits, L. P., editors, *Landscape genetics: concepts, methods, applications*, chapter 14, pages 247–255. Wiley Blackwell, Hoboken, New Jersey.

Balkenhol, N., Gugerli, F., Cushman, S. A., Waits, L. P., Coulon, A., Arntzen, J. W., Holderegger, R., Wagner, H. H., and Participants of the Landscape Genetics Research Agenda Workshop 2007 (2009). Identifying future research needs in landscape genetics: where to from here? *Landscape Ecology*, 24:455–463.

Bishop, C. M. (2006). *Pattern recognition and machine learning*. Springer, New York, New York.

Blei, D. M. (2012). Probabilistic topic models. *Communications of the ACM*, 55:77–84.

Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022.

Carpenter, B., Lee, D., Gelman, A., Goodrich, B., Guo, J., Hoffman, M., Betancourt, M., Li, P., Brubaker, M. A., and Riddell, A. (2015). Stan: a probabilistic programming language. *Journal of Statistical Software*.

Chen, C., Durand, E., Forbes, F., and François, O. (2007). Bayesian clustering algorithms ascertaining spatial population structure: a new computer program and a comparison study. *Molecular Ecology Notes*, 7:747–756.

Coop, G., Witonsky, D., Di Rienzo, A., and Pritchard, J. K. (2010). Using environmental correlations to identify loci underlying local adaptation. *Genetics*, 185:1411–1423.

Cozman, F. G. (2001). JavaBayes: Bayesian networks in Java. `http://www.cs.cmu.edu/~javabayes/index.html`.

Cushman, S. A., McRae, B. H., and McGarigal, K. (2016). Basics of landscape ecology: an introduction to landscapes and population processes for landscape genetics. In *Landscape Genetics: Concepts, Methods, Applications*, chapter 2, pages 11–34. Wiley Blackwell, first edition.

de Guzman, J. and Kaiser, H. (2016). The Boost Spirit library. `http://www.boost.org/doc/libs/1_61_0/libs/spirit/doc/html/index.html`.

de Valpine, P., Turek, D., Paciorek, C., Temple Lang, D., and Bodik, R. (2016). Nimble: numerical inference for hierarchical models using Bayesian and likelihood estimation. `https://bids.berkeley.edu/research/nimble-numerical-inference-hierarchical-models-using-bayesian-and-likeliho`

Dodier, R. (2012). RISO: distributed belief networks. `https://sourceforge.net/`

289     projects/riso/.

290  Dyer, R. J. (2007). The evolution of genetic topologies. *Theoretical Population Biology*, 71:71–79.

291  Dyer, R. J. and Nason, J. D. (2004). Population graphs: the graph theoretic shape of genetic structure.
292     *Molecular Ecology*, 13:1713–1727.

293  Excoffier, L., Smouse, P. E., and Quattro, J. M. (1992). Analysis of molecular variance inferred from
294     metric distances among DNA haplotypes: application to human mitochondrial DNA restriction
295     data. *Genetics*, 131:479–491.

296  Falush, D., Stephens, M., and Pritchard, J. K. (2003). Inference of population structure using
297     multilocus genotype data: linked loci and correlated allele frequencies. *Genetics*, 164:1567–1587.

298  Frichot, E., Schoville, S. D., Bouchard, G., and François, O. (2013). Testing for associations
299     between loci and environmental gradients using latent factor mixed models. *Molecular Biology
300     and Evolution*, 30:1687–1699.

301  Gelman, A., Lee, D., and Guo, J. (2015). Stan: A probabilistic programming language for Bayesian
302     inference and optimization.

303  Gould, S. (2015). DARWIN: a framework for machine learning and computer vision research and
304     development. http://drwn.anu.edu.au.

305  Guillot, G., Estoup, A., Mortier, F., and Cosson, J. F. (2005a). A spatial statistical model for
306     landscape genetics. *Genetics*, 170:1261–1280.

307  Guillot, G., Leblois, R., Coulon, A., and Frantz, A. C. (2009). Statistical methods in spatial genetics.
308     *Molecular Ecology*, 18:4734–4756.

309  Guillot, G., Mortier, F., and Estoup, A. (2005b). GENELAND: a computer package for landscape
310     genetics. *Molecular Ecology Notes*, 5:712–715.

311  Holderegger, R., Kamm, U., and Gugerli, F. (2006). Adaptive vs. neutral genetic diversity: implica-
312     tions for landscape genetics. *Landscape Ecology*, 21:797–807.

313  Holderegger, R. and Wagner, H. H. (2008). Landscape genetics. *BioScience*, 58:199–207.

314  Jia, Y., Salzmann, M., and Darrell, T. (2011). Learning cross-modality similarity for multinomial
315     data. In *Proceedings of the 2011 International Conference on Computer Vision*, ICCV '11, pages
316     2407–2414, Washington, DC. IEEE Computer Society.

317  Joost, S., Bonin, A., Bruford, M. W., Després, L., Conord, C., Erhardt, G., and Taberlet, P. (2007).
318     A spatial analysis method (SAM) to detect candidate loci for selection: towards a landscape
319     genomics approach to adaptation. *Molecular Ecology*, 16:3955–3969.

320  Koller, D. and Friedman, N. (2009). *Probabilistic graphical models: principles and techniques*.
321     MIT Press, Cambridge, MA.

322  Lunn, D., Spiegelhalter, D., Thomas, A., and Best, N. (2009). The BUGS project: evolution, critique
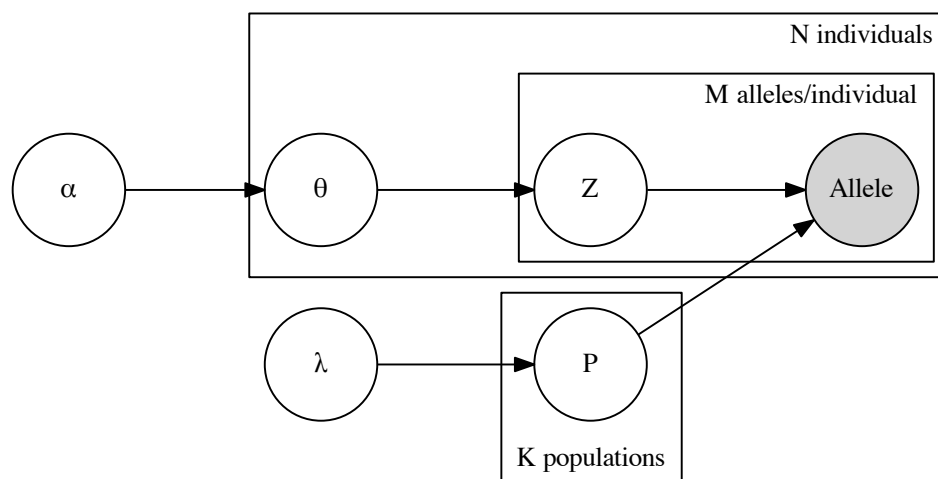323     and future directions. *Statistics in Medicine*, 28:3049–3067.

324  Manel, S. and Holderegger, R. (2013). Ten years of landscape genetics. *Trends in Ecology and
325     Evolution*, 28:614–621.

326  Manel, S., Joost, S., Epperson, B. K., Holderegger, R., Storfer, A., Rosenberg, M. S., Scribner,
327     K. T., Bonin, A., and Fortin, M.-J. (2010). Perspectives on the use of landscape genetics to detect
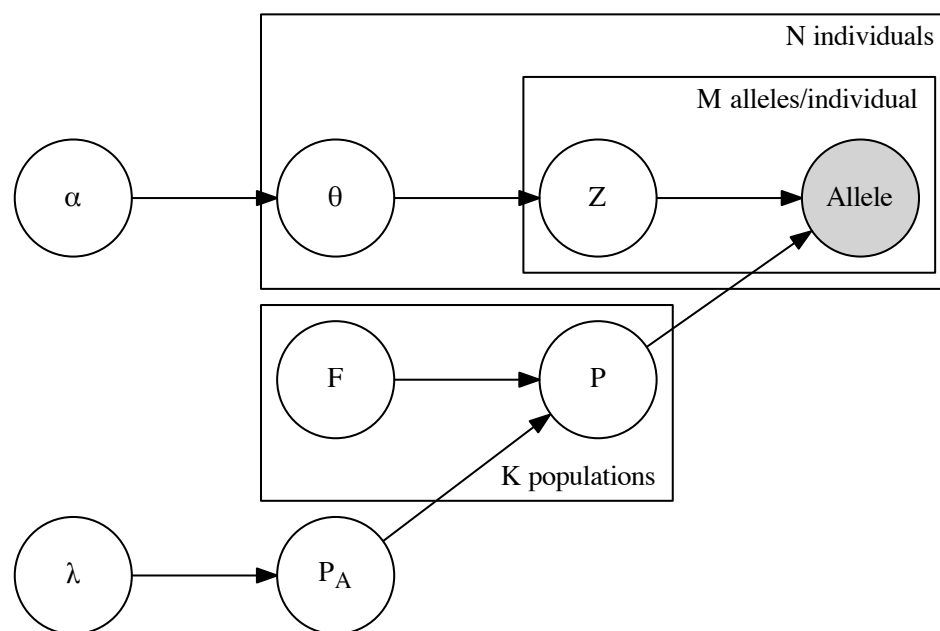328     genetic adaptive variation in the field. *Molecular Ecology*, 19:3760–3772.

329  Manel, S., Schwartz, M. K., Luikart, G., and Taberlet, P. (2003). Landscape genetics: combining
330     landscape ecology and population genetics. *Trends in Ecology and Evolution*, 18:189–197.

331  Minka, T., Winn, J., Guiver, J., Webster, S., Zaykov, Y., Yangel, B., Spengler, A., and Bronskill, J.
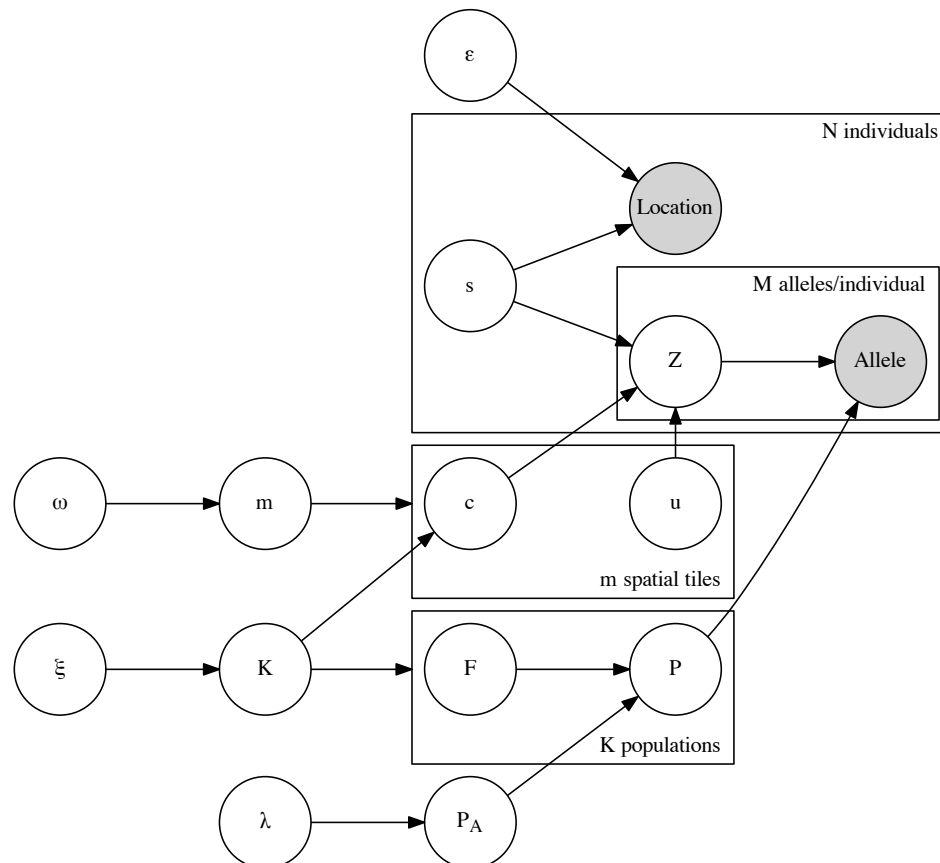
332 (2014). Infer.NET 2.6. `http://research.microsoft.com/en-us/um/cambridge/`
333 `projects/infernet/default.aspx`.

334 Mooiji, J. (2015). libDAI: a free and open source C++ library for discrete approximate inference in
335 graphical models. `https://staff.fnwi.uva.nl/j.m.mooij/libdai/`.

336 Murphy, K. (2014). Software packages for graphical models. `https://www.cs.ubc.ca/`
337 `~murphyk/Software/bnsoft.html`.

338 Murphy, M., Dyer, R., and Cushman, S. A. (2016). Graph theory and network models in landscape
339 genetics. In Balkenhol, N., Cushman, S. A., Storfer, A. T., and Waits, L. P., editors, *Landscape*
340 *genetics: concepts, methods, applications*, chapter 10, pages 165–179. Wiley Blackwell, Hoboken,
341 New Jersey.

342 Niebler, E. (2016). The Boost Proto library. `http://www.boost.org/doc/libs/1_61_`
343 `0/doc/html/proto.html`.

344 OpenGM (2015). OpenGM. `http://hciweb2.iwr.uni-heidelberg.de/opengm/`.

345 Plummer, M. (2015). JAGS version 4.0.0 user manual. Technical report.

346 Plummer, M. (2016). JAGS. `http://mcmc-jags.sourceforge.net`.

347 Pritchard, J. K., Stephens, M., and Donnelly, P. (2000). Inference of population structure using
348 multilocus genotype data. *Genetics*, 155:945–959.

349 Selkoe, K. A., Scribner, K. T., and Galindo, H. M. (2016). Waterscape genetics—applications of
350 landscape genetics to rivers, lakes, and seas. In Balkenhol, N., Cushman, S. A., Storfer, A. T.,
351 and Waits, L. P., editors, *Landscape genetics: concepts, methods, applications*, chapter 13, pages
352 220–246. Wiley Blackwell, Hoboken, New Jersey.

353 Stan Development Team (2016). The Stan math library, version 2.10.0. `http://mc-stan.org`.

354 Stepanov, A. A. and Rose, D. E. (2014). *From Mathematics to Generic Programming*. Addison-
355 Wesley, first edition.

356 Steyvers, M. and Griffiths, T. (2007). Probabilistic topic models. In Landauer, T., McNamara, D.,
357 Dennis, S., and Kintsch, W., editors, *Latent Semantic Analysis: A road to meaning*. Laurence
358 Erlbaum.

359 Storfer, A., Murphy, M. A., Evans, J. S., Goldberg, C. S., Robinson, S., Spear, S. F., Dezzani, R.,
360 Delmelle, E., Vierling, L., and Waits, L. P. (2007). Putting the 'landscape' in landscape genetics.
361 *Heredity*, 98:128–142.

362 Sysoyev, A. V., Milch, B., Bradski, G. R., and Dash, D. (2013). Probabilistic networks library.
363 `https://sourceforge.net/projects/openpnl/`.

364 Thomas, A. (2009). OpenBugs. `http://www.openbugs.net/w/FrontPage`.

365 van Strien, M. J., Keller, D., Holderegger, R., Ghazoul, J., Kienast, F., and Bolliger, J. (2014).
366 Landscape genetics as a tool for conservation planning: predicting the effects of landscape change
367 on gene flow. *Ecological Applications*, 24:327–339.

368 Veldhuizen, T. (1995). Expression templates. *C++ Report*, 7(5):26–31.

369 Warmes, G. (2013). HYDRA MCMC library. `https://sourceforge.net/projects/`
370 `hydra-mcmc/`.

371 Winn, J. (2004). Variational inference for Bayesian networks. `http://vibes.sourceforge.`
372 `net`.

Preprints



**Figure 1.** Plate notation (Bishop, 2006) for the locus-specific graph model used by STRUCTURE (Pritchard et al., 2000). Each circle represents a random variable (or a set of them for those enclosed within boxes) and each arrow represents a dependency of one random variable upon another. This models $N$ individuals each sampled for $M$ (usually two) alleles. $P$ represents the allele frequency distribution in each of $K$ populations and $Z$ represents the assignment of alleles to populations. $\theta$ is the distribution of assignments and $\alpha$ and $\lambda$ are Bayesian priors. The single filled circle indicates that among these random variables only the alleles have been observed; the rest are inferred (or fixed in the case of $\alpha$ and $\lambda$).

**Figure 2.** Plate notation for the correlated allele frequency extension (Falush et al., 2003) to the locus-specific graph model used by STRUCTURE. This models an ancestral population ($P_A$) from which a correlated set of extant populations ($P$) have been derived. The pattern of correlation between populations is governed by $F$.

**Figure 3.** Plate notation for the spatially-explicit extension of STRUCTURE used by GENELAND (Guillot et al., 2005a,b). Additional random variables include the true ($s$) and observed (shaded) locations of sampled individuals and the error ($\varepsilon$) between them, and the locations of points defining the Voronoi tessellation ($u$) and their population identity ($c$). In this case, both the number of Voronoi cells ($m$) and the number of populations ($K$) are random variables.

```
observed_allele_type X;
allele_assignment_type Z;
individual_admixture_distribution_type theta;
population_allele_frequency_distribution_type P;
diriclet_parameter_type alpha;
diriclet_parameter_type lambda;

allele_frequency_type Pr;

for (auto population : populations)
  P(population) =~ dirichlet(lambda);
for (auto individual : individuals)
  {
    theta(individual) =~ dirichlet(alpha);
    for (auto allele : alleles(individual))
      {
        Z(individual,allele) =~ multinomial(theta(individual));
        for (auto population : populations)
          Pr(individual) += Z(population,individual) * P(population);
        X(individual,allele) =~ bernoulli(Pr(individual,allele));
      }
  }
```

**Figure 4.** Compact implementation of the STRUCTURE model with admixture (Pritchard et al., 2000). This is C++ source code for the probabilistic graph model corrresponding to one of the models in STRUCTURE. A few additional lines of code transforms this into a model with correlated allele frequencies (Falush et al., 2003) or one with spatially explicit observations (Guillot et al., 2005a).