

Data-intensive modeling of forest dynamics

Jean F. Liénard^a, Dominique Gravel^b, Nikolay S. Strigul^{a,*}

^a*Department of Mathematics, Washington State University Vancouver, Washington, USA*

^b*Département de Biologie, Université du Québec à Rimouski, Québec, Canada*

Abstract

Forest dynamics are highly dimensional phenomena that are not fully understood theoretically. Forest inventory datasets offer unprecedented opportunities to model these dynamics, but they are analytically challenging due to high dimensionality and sampling irregularities across years. We develop a data-intensive methodology for predicting forest stand dynamics using such datasets. Our methodology involves the following steps: 1) computing stand level characteristics from individual tree measurements, 2) reducing the characteristic dimensionality through analyses of their correlations, 3) parameterizing transition matrices for each uncorrelated dimension using Gibbs sampling, and 4) deriving predictions of forest developments at different timescales. Applying our methodology to a forest inventory database from Quebec, Canada, we discovered that four uncorrelated dimensions were required to describe the stand structure: the biomass, biodiversity, shade tolerance index and stand age. We were able to successfully estimate transition matrices for each of these dimensions. The model predicted substantial short-term increases in biomass and longer-term increases in the average age of trees, biodiversity, and shade intolerant species. Using highly dimensional and irregularly sampled forest inventory data, our original data-intensive methodology provides both descriptions of the short-term dynamics as well as predictions of forest development on a longer timescale. This method can be applied in other contexts such as conservation and silviculture, and can be delivered as an efficient tool for sustainable forest management.

Keywords: data-intensive model, forest dynamics, Gibbs sampling, Markov chain model, Markov chain Monte Carlo, patch-mosaic concept, plant population and community dynamics

*corresponding author

Email address: `nick.strigul@vancouver.wsu.edu` (Nikolay S. Strigul)

1 **Software and data availability**

2 The software to estimate transition matrices based on forest inventory was implemented
3 by Jean Liénard in R version 2.15.1 (R Core Team, 2012) and is attached as a zip file to the
4 submission.

5 The database studied in this paper is available upon request to the Quebec provincial for-
6 est inventory database (<http://www.mffp.gouv.qc.ca/forets/inventaire/>). Straight-
7 forward modifications of the software allows to use with the USDA Forest Inventory and
8 Analysis program (<http://www.fia.fs.fed.us/>).

9 1. Introduction

10 Forest ecosystems are complex adaptive systems with hierarchical structures resulting
11 from self-organization in multiple dimensions simultaneously (Levin, 1999). The patch-
12 mosaic concept was actively developed in the second half of the twentieth century after
13 Watt (1947) suggested that ecological systems can be considered a collection of patches at
14 different successional stages. Dynamical equilibria arise at the level of the mosaic of patches
15 rather than at the level of one patch. The classic patch-mosaic methodology assumes that
16 patch dynamics can be represented by changes in macroscopic variables characterizing the
17 state of the patch as a function of time (Levin and Paine, 1974). Forest disturbances are
18 traditionally associated with a loss of biomass; however, Markov chain models based only
19 on biomass do not capture forest succession comprehensively (Strigul et al., 2012). This
20 limitation motivates the need for alternative formulations that are able to consider several
21 forest dimensions instead of only one.

22 Here we develop a novel statistical methodology for estimating transition probability ma-
23 trices from forest inventory data and generalize classic patch-mosaic framework to multiple
24 uncorrelated dimensions. In particular, we develop a landscape-scale patch-mosaic model
25 of forest stand dynamics using a Markov chain framework, and validate the model using
26 the Quebec provincial forest inventory data. This inventory (Perron et al., 2011) is one of
27 the extensive forest inventories that have been established in North America, among others
28 led by the Canadian provincial governments and the USDA Forest Inventories and Analysis
29 program in the USA. These inventories provide a representative sample of vegetation across
30 the landscape through a large number of permanent plots that are measured repeatedly.
31 Although they were originally developed for estimating growth and yield, they were rapidly
32 found to be extremely useful to studies in forest ecology, biogeography and landscape dynam-
33 ics. Each permanent plot consists of individually marked trees that are periodically surveyed
34 and remeasured. Each plot can be considered as a forest stand and then, theoretically, the
35 forest inventories provide empirical data sufficient for parametrization and validations of
36 patch-mosaic models (Strigul et al., 2012). However, practical development of the patch-

37 mosaic forest models (i.e. their parametrization, validation and prediction) is challenging
38 due to the underlying structure of the forest inventory datasets. These datasets are indeed
39 collected at irregular time intervals that are not synchronized across the focal area, and data
40 collection procedures including spatial plot design and tree measurement methods can be
41 different at various survey times and conducted by different surveyors (Strigul et al., 2012).

42 Our objective in this study is to develop a data-intensive method predicting the dynamics
43 of forest macroscopic characteristics. The idea of a data-intensive modeling approach is to
44 develop and explore a quantitative theory using statistical modeling, in contrast with the
45 hypothesis-driven theoretical approach in which selected mechanisms are used to design
46 and constrain models. We focus here on the development of the modeling framework and
47 illustrate the application of the framework to a large forest inventory dataset spanning 38
48 years of observations collected in Quebec. We rely on Markov chain to describe probabilistic
49 transitions between different states while making minimal assumptions. To overcome the
50 issue of irregular samplings in time specific to forest inventory data, we develop a Gibbs
51 sampling procedure for augmenting the data and infer the transition probabilities. We
52 present in this paper the general methodology and demonstrate each of its steps on the
53 Quebec dataset. In particular, we consider the dimensionality of stand characteristics in this
54 dataset and present evidence that some characteristics are redundant. We apply the method
55 to predict long-term dynamics of Quebec forests, as represented by a subset of macroscopic
56 properties that best represent the variability in the data, and we validate the model utilizing
57 two independent subsets of the original data. We finally discuss the implications of this
58 work, such as the effect of spatial and temporal variability, the independence of most forest
59 variables, the effect of changing external drivers and of feedbacks.

60 **2. Patch-mosaic modeling framework**

61 The goal of this section is to introduce the modeling of patch-mosaic using Markov chains,
62 which is generalized and employed to predict forest dynamics in the main text. The patch-
63 mosaic concept assumes that the vegetation at the landscape level can be represented as

64 a collection of isolated spatial units - patches - where patch development follows a general
65 trajectory and is subject to disturbances (Watt, 1947; Levin and Paine, 1974). Patch-mosaic
66 models are derived using the conservation law, which takes into account patch aging and
67 other changes to macroscopic variables representing succession, growth of patches in space,
68 and disturbances (Levin and Paine, 1974). The same general idea as well as mathematical
69 derivations are broadly used in population dynamics to describe age- and size-structured
70 population dynamics. Patch-mosaic models can be partial differential equations or discrete
71 models depending on whether time and patch stages are assumed to be continuous or discrete.
72 Classic continuous patch-mosaic models are based on the application of the conservation
73 law to continuously evolving patches that can be destroyed with a certain probability, and
74 can be represented by the advection equation (model developed by Levin and Paine, 1974,
75 for fixed-size patches) or equivalently by the Lotka-McKendrick-von Foerster model (Strigul
76 et al., 2008). The continuous patch-mosaic models have been used in forest ecology to model
77 the dynamics of individual canopy trees within the stand or forest gap dynamics (Kohyama
78 et al., 2001; Kohyama, 2006).

79 In the case of patches changing in discrete time, the derivation of the conservation law
80 leads to discrete-type patch-mosaic models. In particular, the advection-equation model
81 (Levin and Paine, 1974) is essentially equivalent to several independently developed discrete
82 models (Leslie, 1945; Feller, 1971; Van Wagner, 1978; Caswell, 2001). These models consider
83 only large scale catastrophic disturbances (patch "death" process), destroying the patch,
84 which then develops along the selected physiological axis until the next catastrophic dis-
85 turbance Levin and Paine (1974). The stochastic model we are considering here employs a
86 Markov chain framework (Waggoner and Stephens, 1970; Usher, 1979a; Facelli and Pickett,
87 1990; Logofet and Lesnaya, 2000; Caswell, 2001) that is capable of taking into account all
88 possible disturbances.

89 In a Markov chains model, the next state of a forest stand depends only on the previous
90 state, and the probabilities of going from one state into another are summarized in what is
91 called a transition matrix, denoted T .

92 We summarize the distribution of states at time t as the row vector X_t , with length equal
93 to the number of discrete classes of patch state and with a sum equal to 1. We can predict
94 $X_{t+\Delta t}$ by multiplying the transition matrix:

$$X_{t+\Delta t} = X_t \cdot T \quad (1)$$

95 To project an arbitrary number n time steps into the future, one simply multiplies by
96 T^n instead of T . The Perron-Frobenius Theorem guarantees the existence of the long-term
97 equilibrium, which can be practically found as the normalized eigenvector corresponding to
98 the first eigenvalue, or by iterative sequence of state vectors. In this paper we employ the
99 iterative method as it allows to derive forest states at different time steps in the future, for
100 example allowing to make predictions in 10, 20 or 30 years from now. To derive the long-term
101 equilibrium we simply choose an n large enough to satisfy the condition:

$$|X_{t+n\Delta t} - X_{t+(n-1)\Delta t}| < \epsilon \quad (2)$$

102 Three illustrative examples of simplified Markov chains are available in Appendix 1.

103 3. Materials and Methods

104 We address here the issue of constructing transition matrices from forest inventory
105 data stemming from irregular sampling intervals and variable numbers of plots sampled in
106 each year. We outline in the following the general concepts of the methodology along with
107 practical guidelines using the inventory led by the provincial Ministry of Natural Resources
108 and Wildlife in Quebec (Appendix 1). The key steps to use Gibbs sampling to estimate a
109 transition matrix from irregular measurements are:

- 110 1. Compute stand level characteristics for each plot and for each survey year. Analyze
111 the dimensionality of these characteristics using correlation and principal component
112 analysis;

113 2. Construct temporal sequences of uncorrelated characteristics depending on forest sur-
114 vey dates. Use Gibbs sampling to infer the transition matrix. This algorithm consists
115 of random initialization of missing values followed by iteration of parameter estimation
116 and data augmentation:

- 117 • Parameter estimation: Compute the transition matrix using the (augmented)
118 sequences of plot transitions.
- 119 • Data augmentation: Draw new sequences conditional on the new transition ma-
120 trix.

121 The transition matrices for Quebec forests were obtained using this method with a three-
122 year time step. Future and equilibrium landscape characteristics were predicted according
123 to equations 1 and 2 (cf. Appendix 1).

124 *3.1. Step 1: stand characteristics and dimensional analysis*

125 This step consists of (a) the selection of a set of stand-level forest characteristics, (b) the
126 dimensional analysis of these characteristics, (c) their decomposition into uncorrelated axes,
127 and (d) the discretization of these uncorrelated axes.

128 Our modeling method can be applied for the prediction of any forest stand characteristic
129 under the condition that it is computable from every single plot survey. The particular choice
130 of the characteristics depends on available data and research objectives. A general guideline
131 is that these characteristics should summarize data from individual trees into macroscopic
132 indicators of stand structure, which can then be used to compare forests across different
133 ecosystems. We consider six characteristics of Quebec forests according to the rationale pre-
134 sented in Strigul et al. (2012) and Lienard et al. (2014). These characteristics are computed
135 based on trees with a diameter at breast height larger than 90mm (see Appendix 1.1 for
136 more details about the Quebec forest inventory measuring protocol). We denote \mathcal{S} the set of
137 species inside each plot and \mathcal{T} the set of trees inside each plot, and compute for each single
138 plot survey the following characteristics:

- 139 • dry weight biomass, estimated from Jenkins et al. (2003), using the formula: $\sum_{i \in \mathbb{T}} e^{B1_i + B2_i \log(d_i)}$
140 where B1 and B2 are species specific density constants, and d is the trunk diameter
141 at breast height in cm. B1 and B2 have been derived from both US and Canadian
142 studies, making it a suitable approximation for Quebec forests (Jenkins et al., 2003).
143 The resulting aboveground biomass is expressed in 10^3 kg/ha.
- 144 • basal area, computed as the sums of trunk diameters at breast height d : $\sum_{i \in \mathbb{T}} \pi \left(\frac{d_i}{2}\right)^2$.
145 The basal area is expressed in m^2/ha .
- 146 • intra-plot diversity (evenness), computed as the Gini-Simpson index (Hill, 2003), with
147 $\Omega(s)$ referring to the number of trees with species s and $\Omega(\mathbb{T})$ referring to the total
148 number of trees inside each plot: $1 - \sum_{s \in \mathbb{S}} \left(\frac{\Omega(s)}{\Omega(\mathbb{T})}\right)^2$. This provides an index in the
149 0-1 range describing the species heterogeneity at the stand level, with high values
150 indicating a high heterogeneity.
- 151 • extra-plot diversity (species richness), computed as the number of species present in a
152 plot: $\Omega(\mathbb{S})$. In the Quebec dataset, this indicator ranges from 1 to 8 species, and is
153 interpreted as another measure of diversity.
- 154 • shade tolerance index, a new metric introduced by Strigul and Florescu (2012) and
155 Lienard et al. (2014) describing the shade tolerance rank of species r : $\sum_{i \in \mathbb{T}} \frac{\Omega(s)r_i}{\Omega(\mathbb{T})}$.
156 This index ranges from 0 to 1, with high values denoting forest stands composed of
157 typically late successional species and low values denoting forest stands composed of
158 typically early successional species in Quebec (Lienard et al., 2014).
- 159 • average age, computed as the average of tree ages a : $\sum_{i \in \mathbb{T}} \frac{a_i}{\Omega(\mathbb{T})}$. This commonly-used
160 indicator approximates the stand age in the forest inventory analysis (see Strigul et al.
161 2012 for a discussion of this characteristic).

162 Statistical relations of these stand-level characteristics were analyzed using standard mul-
163 tivariate methods. First, we computed the Pearson correlation coefficients both in the whole
164 dataset and in the dataset broken down in decades (to avoid biases due to their temporal

165 autocorrelation). We then performed a principal component analysis (PCA) to examine
166 (a) the number of components needed to explain most of the variance as well as (b) the
167 projection of characteristics in the space defined by these components.

168 In general, it is possible for a multidimensional model to operate on the space of principal
169 components. Such a model would (a) project the characteristics into the low-dimensional
170 space given by the principal components, then (b) predict their dynamics in this new space,
171 and finally (c) perform the inverse transformation to obtain predictions on the characteristics.
172 In our application to the Quebec dataset, we discovered that four uncorrelated characteristics
173 approximate well the principal component space (namely biomass, average age of trees, Gini-
174 Simpson and shade tolerance indexes, *cf.* Results). Our model employs this approximation
175 and is based on transition matrices of these forest characteristics. It substantially simplify
176 interpretation of modeling predictions.

177 Prior to the computation of transition matrices in the Markov chain framework, it is
178 necessary to discretize continuous variables into distinct states (Strigul et al., 2012). The
179 general approach is to subdivide data into uniformly spaced states, with a precision that is
180 small enough to capture the details of the distribution but large enough to be insensitive
181 to statistical noise in the dataset. In addition, the computational effort needed to infer
182 transition matrices is proportional to the square of the number of states, and available
183 computational power may constitute a practical limitation to the number of states. In the
184 Quebec dataset, the stand-level characteristics span different ranges (see Figs. 2 and 3 in
185 Appendix), with the biomass distribution in particular showing a long tail for the highest
186 values. In order to capture enough details of the distributions of the Quebec characteristics,
187 we opted to remove plots in the long tail of the biomass (those with a biomass higher
188 than 50,000 kg/ha, representing roughly 4% of the total dataset) and then subdivided the
189 remaining plots into 25 biomass states. An alternative approach would be to merge the rarely
190 occurring high-biomass states into the last state as was implemented in Strigul et al. (2012).
191 We conducted a comparison of these two approaches and found no significant differences.
192 For the other characteristics investigated (i.e. the internal diversity, shade tolerance index,

193 and average age), we found that 10 states were enough to capture their distributions with
194 sufficient detail.

195 *3.2. Step 2: Gibbs sampling methodology*

196 Inferring a Markov Chain model for characteristics computed with field data sampled at
197 irregular intervals is a challenging problem. Indeed, the usual direct approach of establishing
198 the n -year transition matrix by simply counting the number of times each state changes to
199 another after n years can not be employed in most forest inventories, as successive mea-
200 surements on the same plot are not made with constant time intervals. This irregularity
201 in sampling results in states of the forest plots that are not observed, and can be modeled
202 as missing data. Two classes of algorithms can be used to parameterize a transition matrix
203 describing the dynamics of both observed and missing data: expectation-maximization (EM)
204 and Monte Carlo Markov Chain (MCMC), of which Gibbs sampling is a specific implemen-
205 tation. Both classes of algorithms are iterative and can be used to find the transition matrix
206 that best fits the observed data. EM algorithms consist of the iteration of two steps: in the
207 expectation step the likelihood of transition matrices is explicitly computed given the distri-
208 bution of the missing data inferred from the previous transition matrix estimate, and in the
209 maximization step a new transition matrix maximizing this likelihood is chosen as the new
210 estimate (Dempster et al., 1977). MCMC algorithms can be seen as the Bayesian counter-
211 part of EM algorithms, as at each iteration a new transition matrix is stochastically drawn
212 with the prior information of estimated missing data, and in turn new estimates for the
213 missing data are stochastically drawn from the new transition matrix (Gelfand and Smith,
214 1990). EMs are deterministic algorithms, and as such they will always converge to the same
215 transition matrix with the same starting conditions; conversely, MCMCs are stochastic and
216 are not guaranteed to converge toward the same estimate with different random seeds. While
217 both algorithms are arguably usable in our context, the ease of implementation and lower
218 computational cost of MCMC algorithms led us to prefer them over EM (Deltour et al.,
219 1999). We selected Gibbs sampling as a flexible MCMC implementation (Geman and Ge-
220 man, 1984). We provide in the following a brief presentation of Gibbs sampling. Additional

221 implementation details are in Appendix 1.2, and we refer to Robert and Casella (2004) for
222 the general principles underlying MCMC algorithms and to Pasanisi et al. (2012) for an ex-
223 tended description of Gibbs sampling to infer transition probabilities in temporal sequences.
224 In addition to the full explanation below, we also provide a pseudocode of the procedure
225 (Box 1).

226 To apply Gibbs sampling for the estimation of the transition matrices, it is required to
227 include plot characteristics in a set of temporal sequences. For each plot p , this is done by
228 inserting each characteristic $s_{(p,i)}$ measured in the i -th year at position i of a row vector S_p
229 representing the temporal sequence of this plot. For example, if a plot p was sampled only
230 at years 1 and 3 during a 5-year inventory, allowing for the computation of characteristics
231 $s_{(p,1)}$ and $s_{(p,3)}$, then its sequence would be the row vector $S_p = [s_{(p,1)}, \bullet, s_{(p,3)}, \bullet, \bullet]$, where
232 \bullet denotes a missing value. The sequences are mostly composed of unknown values as only
233 a fraction of the forest plots were surveyed each year. In the application to the Quebec
234 dataset, a reduction of the size of these temporal sequences was performed (see Appendix
235 1.2 for a detailed description of this reduction and an illustrative example); however it is not
236 a pre-requisite for the general application of Gibbs sampling. Let further Y be the matrix
237 constructed using all the sequences S , with rows corresponding to successive measures of
238 different plots and columns corresponding to different years. The preliminary step of Gibbs
239 sampling consists of replacing the missing values \bullet in Y at random, resulting in so-called
240 augmented data $Z^{[0]}$. Then, the two following steps are iterated a fixed number of times H ,
241 with h the index of the current iteration:

- 242 1. in the **parameter estimation** step, we draw a new transition matrix $\Phi^{[h]}$ conditional
243 on the augmented data Z^{h-1} , using for every row i :

$$\Phi_i^{[h]} | Z^{[h-1]} \sim Dir(\gamma_{i,1} + w_{i,1}^{[h-1]}, \dots, \gamma_{i,r} + w_{i,n}^{[h-1]}) \quad (3)$$

244 with Dir is the Dirichlet distribution, γ are biasing factors set here uniformly to 1 as
245 we include no prior knowledge on the shape of the transition matrix (Pasanisi et al.,
246 2012). $w_{i,j}$ are the sufficient statistics reflecting the transitions in the augmented data

247 $Z^{[h-1]}$, formally defined as

$$w_{i,j} = \sum_{t \in \text{years}} \sum_{k \in \text{plots}} \mathbb{1}_{\{Z_{k,t-1}^{[h-1]}=s_i \ \& \ Z_{k,t}^{[h-1]}=s_j\}} \quad (4)$$

248 with $\mathbb{1}_{\{Y_{k,t-1}=s_i \ \& \ Y_{k,t}=s_j\}}$ the count of sequences elements in the state s_i at time $t - 1$
 249 and in the state s_j at time t .

250 2. in the **data augmentation** step, we draw new values the missing states, based to the
 251 probabilities of the transition matrix $\Phi^{[h]}$. The probabilities \mathbb{P} used to augment the
 252 data $Z^{[h]}$ are derived from their values in the previous iteration ($Z^{[h-1]}$) as well as their
 253 values in the current iteration but in an earlier year ($Z_{k,t-1}^{[h]}$ for $t \geq 2$):

$$\text{for the earliest data } t = 1, \quad \mathbb{P}(Z_{k,1}^{[h]} = s_j | Z_{k,2}^{[h-1]} = s_i, \Phi^{[h]}) \propto \Phi_{j,i}^{[h]} \quad (5)$$

$$\text{for the latest data } t = T, \quad \mathbb{P}(Z_{k,T}^{[h]} = s_j | Z_{k,T-1}^{[h]} = s_i, \Phi^{[h]}) \propto \Phi_{i,j}^{[h]} \quad (6)$$

$$\text{otherwise, } \mathbb{P}(Z_{k,t}^{[h]} = s_j | Z_{k,t-1}^{[h]} = s_{i_1}, Z_{k,t+1}^{[h-1]} = s_{i_2}, \Phi^{[h]}) \propto \Phi_{i_1,j}^{[h]} \times \Phi_{j,i_2}^{[h]} \quad (7)$$

Pseudocode 1: Estimate of the transition matrix of one stand characteristic

Data: Y , a matrix whose rows are sequences S of repeated measurements.

Result: M , the transition matrix

begin

transitions list $\leftarrow \emptyset$;

for $r \leftarrow 1$ **to** R **do**

/ initialization* **/*

 Obtain $Z^{[0]}$ by filling missing states from Y at random;

for $h \leftarrow 1$ **to** H **do**

/ Parameter estimation* **/*

 Compute the sufficient statistics $w_{i,j}, \forall (i, j) \in [1, n]^2$ using Eq. 4;

 Initialize $\Phi^{[h]}$ as an empty $n \times n$ matrix;

for $i \leftarrow 1$ **to** n **do**

$\Phi_{i,1\dots n}^{[h]} \leftarrow$ random values drawn from Eq. 3;

end

/ Data augmentation* **/*

$Z^{[h]} \leftarrow Y$;

for $k \leftarrow 1$ **to** K **do**

if missing, $Z_{k,1}^{[h]} \leftarrow$ random value drawn from Eq. 5;

end

for $t \leftarrow 2$ **to** $T - 1$ **do**

for $k \leftarrow 1$ **to** K **do**

if missing, $Z_{k,t}^{[h]} \leftarrow$ random value drawn from Eq. 6;

end

end

for $k \leftarrow 1$ **to** K **do**

if missing, $Z_{k,T}^{[h]} \leftarrow$ random value drawn from Eq. 7;

end

end

if $h > B$ **then**

 transitions list \leftarrow {transitions list, $\Phi^{[h]}$ };

end

end

end

$M \leftarrow$ mean of all matrices in transitions list

254 As Gibbs sampling is initialized by completing the missing values at random, the first
255 iterations will likely result in transition matrices far away from the optimal. The usual
256 workaround is to ignore the first B transition matrices corresponding to so-called "burn-in"
257 period, leaving only $H - B$ matrices. Furthermore, as Gibbs sampling relies on a stochastic
258 exploration of the search space, a good practice to ensure that Gibbs sampling converged
259 to the optimal transition matrix is to run the whole algorithm R times. There are no
260 general guidelines for setting the H , B and R parameters (Robert and Casella, 2004). We
261 empirically settled with $H = 1000$, $B = 100$ and $R = 50$ in order to ensure that the transition
262 matrices were reproducible for the Quebec dataset, leading to $R \times H = 50000$ iterations of
263 parameter estimation and data augmentation steps and resulting in $R \times (H - B) = 45000$
264 transition matrices. This process was repeated independently for each plot characteristic.
265 The algorithm was implemented in R version 2.15.1 (R Core Team, 2012) and took a total
266 runtime of 4 days on a 1.2 Ghz single-core CPU to compute the transition matrices for all 4
267 characteristics studied here.

268 4. Results

269 4.1. Multivariate analysis of stand characteristics

270 The correlation analysis performed on the Quebec forest inventory (Perron et al., 2011,
271 Appendix 1.1) revealed that biomass and basal area were highly correlated ($r = 0.96$), as
272 well as the external and internal diversity indices ($r = 0.90$, see Appendix 1.3 for the other
273 coefficients). These correlations are further preserved when the correlation analysis is done
274 separately on each decade, from the 1970s until the 2000s (*cf.* tables in Appendix 1.3),
275 confirming the presence of time-independent strong correlations between these two pairs of
276 characteristics.

277 A PCA applied to the dataset further confirmed that the biomass and basal area on one
278 hand, as well as the external and internal diversity on the other hand, have nearly identical
279 vectors in the principal components space (*cf.* Appendix 1.4). Furthermore, this analysis
280 showed that 4 principal components are required to adequately explain variance in the data;

281 using 3 components accounts for only 87 % of the variance, while 4 components explain up
282 to 98 % of the variance. The PCA revealed that biomass, the internal diversity index, the
283 shade tolerance index, and the average age are close approximations of the different principal
284 components and explain most of the variance. Therefore, these variables have been employed
285 in the following analysis.

286 *4.2. Interpretation of the transition matrices*

287 We present here in detail the transition matrix for biomass with a 3-year time interval,
288 shown in Fig. 1 (the other characteristics are to be found in Appendices, in Figs. 7 and 8).
289 In this matrix, each value at row i and column j corresponds to the probability of transition
290 from state i into state j after 3 years. By definition, rows sum to 100%. This transition
291 matrix, as with the others in Appendix, is dominated by its diagonal elements, which is
292 expected because few plots show large changes in a given 3-year period. The values below
293 the diagonal correspond to transitions to a lower state (hence, they can be interpreted as
294 the probabilities of disturbance), while values above the diagonal correspond to transitions
295 to a higher state (i.e., growth). The transitions in the first column of the matrix correspond
296 to major disturbances, where the stand transitions to a very low biomass condition. As
297 the probabilities above the diagonal are larger than below the diagonal, the overall 3-year
298 prediction is of an increase in biomass. This matrix also shows that plots with a biomass
299 larger than 40,000 kg/ha have a roughly uniform 10% probability of ending with a biomass of
300 less than 20 000 kg/ha 3 years later, which is interpreted as the probability of high-biomass
301 stand to go through a moderate to high disturbance.

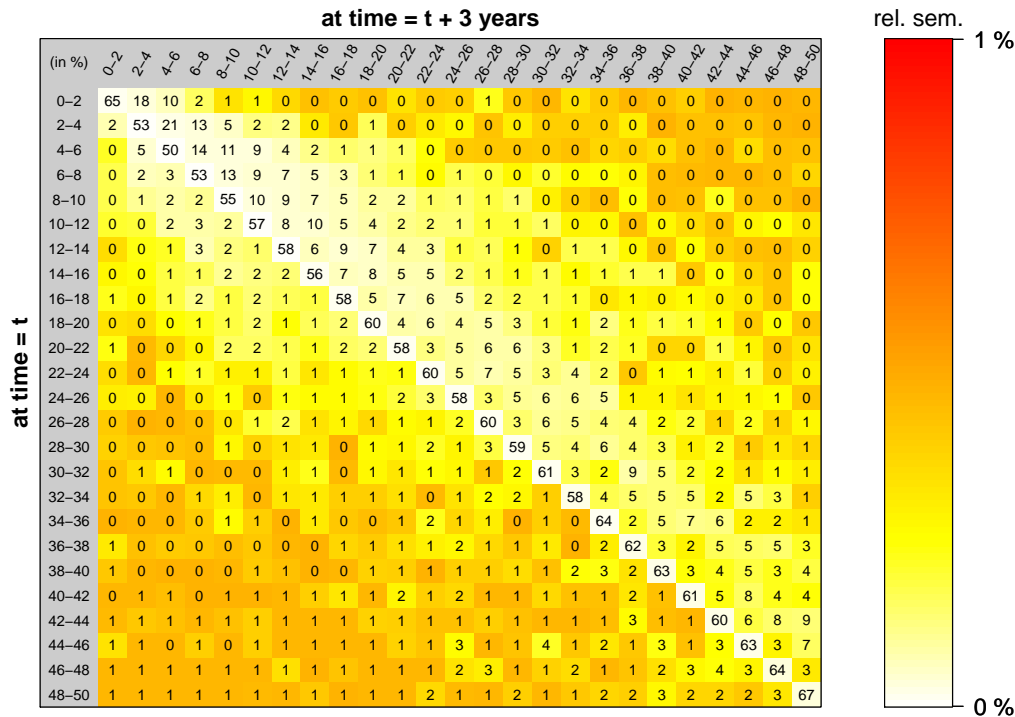


Fig. 1. 3-year transition matrix for the biomass. The states are the biomass ranges in 10^3 kg/ha, spanning from 0 – 2 to 48 – 50 10^3 kg/ha, and represented here on the left and on top of the matrix. The values $M(i, j)$ inside the matrix correspond to the rounded probability of transition from state i to state j . The color represents the relative standard error of the mean and indicates the robustness of the stochastic search, as explained in Section 4.3. Lighter colors thus indicate a better confidence in the transition value; all relative standard errors of the mean (RSEM) are below 1%, corresponding to a very high confidence, and furthermore the smallest errors are found for the higher transition probabilities close to the diagonal.

302 4.3. Model validation

303 Two main types of error should be considered when designing a model with a parameter
 304 search based on real data. The first error relates to the robustness and efficiency of the
 305 estimation of the optimal transition matrix, which was performed with Gibbs sampling in
 306 our case. The second type of error encompasses more broadly the capacity of the chosen
 307 theoretical framework to predict the system beyond the range of the dataset. In our case, the
 308 theoretical framework we relied on is patch-mosaic concept, implemented with the Markov
 309 chain machinery, to describe the dynamics of our four characteristics.

310 To estimate the errors of the parameter search, we used the $R(H - B)$ transition matrices to compute for each transition the standard error of the mean (SEM) and the relative
311 standard error of the mean (RSEM, defined as the ratio of the SEM over the transition probability, and expressed as a percentage). The SEM were below 1% throughout the matrices,
312 with the highest errors occurring for very low transition probabilities (i.e., far from the diagonals). Furthermore, the RSEM were very low, and particularly so for the transitions with
313 the highest probability (Fig. 1 in main text as well as Figs. 7 and 8 in Appendix). We finally
314 computed the SEM in the long-term predicted equilibriums and found values below 0.01%,
315 strengthening the conclusion that negligible errors are to be attributed to the stochastic fit
316 procedure.

320 An independent dataset would be most suited to estimate the more general errors in
321 the ability of a Markov-Chain model to predict future forest characteristics. As there is no
322 such dataset available, we performed two cross-validations of our methodology by splitting
323 this dataset in two different ways. In the first, we ran the Gibbs sampler with only the
324 first 18 years of records (from 1970 to 1988). We then used the model to predict forest
325 state for the period corresponding to the second half of the dataset (i.e., 1989 to 2007), and
326 we compared the predicted dynamics with the aggregated distribution of the second half
327 of the dataset (Fig. 2). Overall, the predictions were highly accurate, with R^2 between
328 observation and prediction ranging from 0.8 to 0.95, indicating that the second half of the
329 dataset is predictable with a Markov chain model based solely on the first half. In the
330 second validation, we randomly split the data into two sets, regardless of year. We then
331 computed the transition matrix and corresponding equilibrium conditions for each half (Fig.
332 9 in Appendix). Here again, the predictions match closely with values of R^2 higher than 0.98
333 for the internal diversity, shade tolerance index and average age. The R^2 was near 0.6 for
334 the biomass, indicating that this variable is more sensitive to small changes than the others;
335 however the difference in predictions were small, typically around 1% for each biomass state.
336 This second validation overall showed that the data contained in the inventory is redundant,
337 and that half of it is enough to provide highly accurate long-term estimates for the internal

338 diversity, shade tolerance index and average age. Considering only half of the data at random
 339 would likely result in errors of around 1% in the long-term estimates of the biomass.

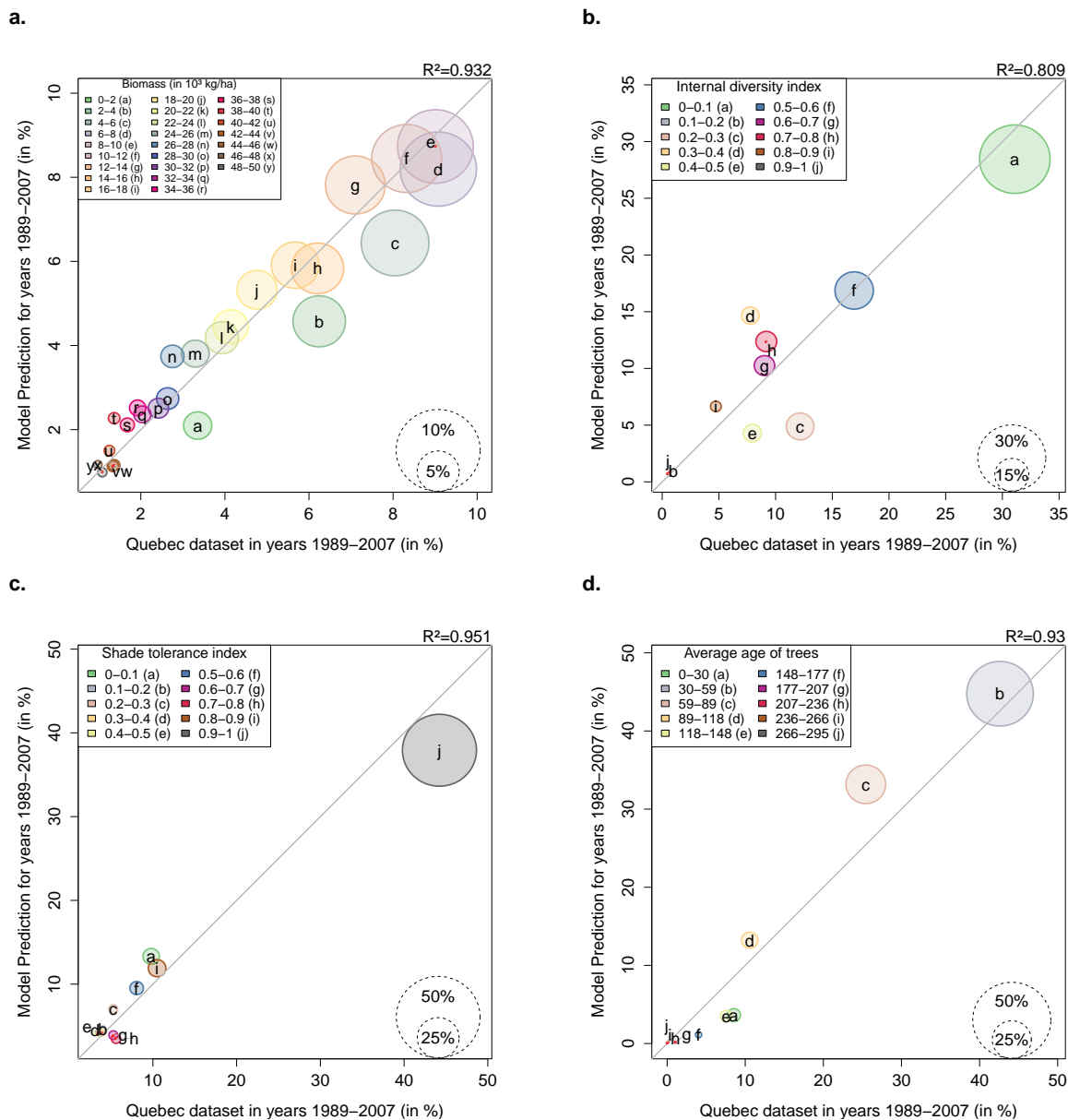


Fig. 2. Results of model validation, showing the second half of the dataset *vs* the model prediction for the classes of each characteristic (distribution in %). For each class, the circle size denotes the number of stands belonging to it in the real dataset. The R^2 measure is indicated on the top right of each plot. The model used to make the prediction was computed using only the first half of the dataset, corresponding to years 1970 to 1988 (see Materials & Methods for details).

340 *4.4. Predictions of temporal dynamics and long-term equilibrium*

341 We applied the inferred transition matrix to predict the state of forest in 2010s, 2020s
342 and 2030s based on their distribution in 2000s. We also predicted the long-term dynamics
343 of the forest stands, by computing the equilibrium states of the transition matrices. Overall,
344 the predictions showed an increase in biomass and stand age (Fig. 3 e and h), along with
345 a slight increase in biodiversity (Fig. 3 f) and a slight decrease of the prevalence of late
346 successional species accompanied by a slight increase of early successional species (Fig. 3 g).
347 These predictions are obvious for the biomass and average age of trees by looking at their
348 distributions in the existing dataset (Fig. 3 a and d), while they are less clearly seen when
349 looking at the average distributions of the biodiversity and shade tolerance index (Fig. 3 b
350 and c).

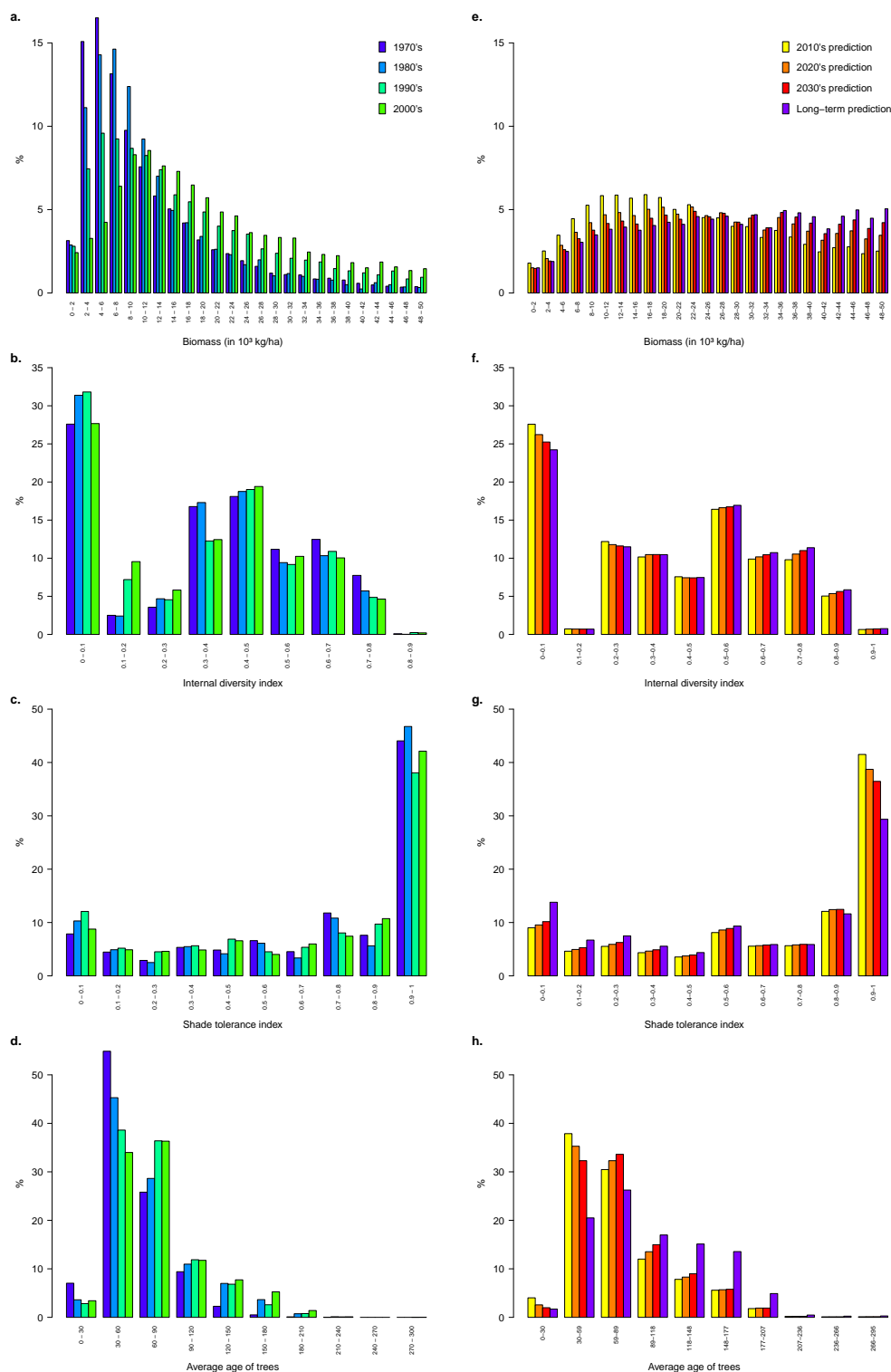


Fig. 3. Current distribution of relevant characteristics from the database, along with the long-term predictions of our models.

351

These long-term predictions were reached at different timescales depending on the char-

acteristics. For biomass, equilibrium was reached by approximately year 2030, but the other characteristics, and in particular the average age of trees in plots, showed much slower dynamics to reach their equilibria (Fig. 3 e to h). The model predicted average relative changes of +38.9% and +14.2% by the 2030s for biomass and stand age, and +44.0% and +37.9% by the time they reach their long-term equilibrium state. Relative changes for the Gini-Simpson diversity index were +5.2% by the 2030s and +7.1% in the long term, and early successional species will become slightly more abundant with a change of -4.7% of the shade tolerance index by the 2030s and -13.6% in the long term.

5. Discussion

We developed a data-intensive approach to multiple-dimensional modeling of forest dynamics. The modeling steps include 1) dimensional analysis of forest inventory data, 2) extraction of non-correlated dimensions, and 3) the application of stochastic optimization to compute probability transition matrices for each dimension. We applied this approach to the Quebec forest inventory dataset and validated the model using two independent subsets of data. Our study demonstrates that there exist at least four uncorrelated dimensions in Quebec forests: the biomass, biodiversity, shade tolerance index and averaged age of trees. The most pronounced changes predicted for Quebec forests are increases in biomass and stand age. Our model also predicted smaller increases in biodiversity in the prevalence of early successional species. Our results demonstrate the utility of this methodology in predicting long-term forest dynamics given highly dimensional, irregularly sampled data; the model was computationally efficient and validation procedures demonstrated its ability to make short and long-term predictions. Therefore, the framework will be useful both in applied contexts (e.g., conservation, silviculture) as well as in developing our conceptual understanding of how forested ecosystems are organized through dimensional analysis of forest characteristics under the current disturbance regime.

377 *5.1. Contribution to Markov chain forest modeling framework*

378 Markov chain models have a rich history of application in ecology, and, in particular,
379 in forest modeling (Facelli and Pickett, 1990; Caswell, 2001). This modeling framework
380 has been employed to describe forest transitions at different scales with various focal vari-
381 ables, for example, succession models defined on the species and forest type level (Usher,
382 1969, 1981; Waggoner and Stephens, 1970; Horn, 1974; Logofet and Lesnaya, 2000; Ko-
383 rotkov et al., 2001), gap mosaic transition models (Acevedo et al., 1996, 2001) and biomass
384 transition models (Strigul et al., 2012). Markov chain successional models (Usher, 1969,
385 1981, 1979b; Facelli and Pickett, 1990) are able to predict changes in species abundance,
386 but require a comprehensive knowledge of successional sequence of species replacement and
387 transition probabilities between different successional stages. The empirically-based Markov
388 chain forest succession model, which operates at the species level and assumes that the un-
389 derlying Markov chain is stationary, requires only substantially large observations to estimate
390 transition probabilities (Waggoner and Stephens, 1970; Stephens and Waggoner, 1980). On
391 the contrary, the mechanistic Markov chain modeling approach developed by Horn (1974,
392 1981) employs shade tolerance and gap dynamics to predict species replacement in the forest
393 canopy given the species composition in the understory. However, this approach requires
394 a detailed survey of the understory vegetation that is not commonly available in forest in-
395 ventories. Also, gap dynamics individual-based models can be coupled with Markov chain
396 models for scaling of gap dynamics to patch level (Acevedo et al., 1996, 2001). These tran-
397 sitional models have been demonstrated to be useful and relevant tools in forest prognosis,
398 however their practical applications are often limited. The Bayesian methodology proposed
399 in this study allows to extend the scope of transition matrices by allowing their computation
400 directly from forest inventory data, with corresponding modifications of the R code provided
401 as a supplementary material. The proposed methodology of matrix estimation could be
402 employed to test the validity of the Markov chain homogeneity assumption.

403 *5.2. A data-intensive approach to understand forest dynamics*

404 Modeling complex adaptive systems such as forest ecosystems requires capturing the
405 dynamics of biological units at multiple scales and in multiple dimensions (Levin, 1998,
406 2003). Ideally, a mechanistic model based on the physiological processes and interactions
407 of individual organisms should simulate the observed forest structure and predict forest
408 dynamics over different time horizons and environmental variables. However, such individual-
409 based modeling is very challenging as interactions between individual organisms within the
410 forest stand result in new properties at the stand level, where essential mechanisms, spatial
411 dimensions of variables, and functional relationships between variables are largely unknown.
412 Given these unknowns, a data-intensive approach can be useful for gaining insight into
413 ecosystem dynamics provided that sufficient amounts of relevant data are available (Kelling
414 et al., 2009; Michener and Jones, 2012). In particular, the matrices
415 we estimated (see Fig. 1 in main text and Figs. 7 and 8 in Appendices) incorporate all
416 forest changes related to different magnitude disturbances. This opens a possibility of the
417 future investigation of how particular disturbances are reflected in the forest macroscopic
418 characteristics and can lead to a logical extension of classic models that take into account
419 only major disturbances, in particular, birth-and-disaster Markov chains (Feller, 1971), forest
420 fire models (Van Wagner, 1978), and advection-reaction equations for patch dynamics (Levin
421 and Paine, 1974) .

422 A potential limitation to a mechanistic interpretation of the transition matrices arises
423 from the Markovian assumption that the transition toward the next state depends solely
424 on the current state. If this assumption is not valid, it could bias these models. This
425 assumption warrants further attention as it has not been yet comprehensively evaluated in
426 forest modeling.

427 Integral Projection Model (IPM) is another modeling framework that could be used in
428 place of Markov chains (Easterling et al., 2000; Caswell, 2001). In IPM, continuous kernel
429 functions are used instead of discrete transition probabilities. While IPM are by design suited
430 to handle well data irregularly distributed across the states, they do not address explicitly

431 the issue of sampling irregularities in time. The data augmentation approach developed here
432 can however be transposed to parameterizing IPM as well. Markov chains are preferable in
433 our application because they are not restricted by the choice of IPM kernels. Indeed, biomass
434 and stand age transitions can be decomposed into several kernels using commonly accepted
435 assumptions of growth and disturbances, however there is no obvious way to choose kernels
436 for biodiversity and shade tolerance index as their dynamics can not be understood in terms
437 of a monotonic progression toward high values (Lienard et al., 2014).

438 The application of MCMC procedures allows to compute transition matrices for datasets
439 with irregular sampling intervals and sample sizes. While Gibbs sampling has been intro-
440 duced 30 years ago (Geman and Geman, 1984), its application to handle missing data in
441 ecology has been mostly limited to stochastic patch occupancy models with a low number of
442 free parameters (5-6) and either artificially simulated data or relatively restricted datasets
443 (e.g. 72-228 resampled locations in ter Braak and Etienne, 2003; Harrison et al., 2011; Risk
444 et al., 2011). From the technical point of view, our application of MCMC differs by taking
445 advantage of the absolute time independence of Markov chains (allowing us to align sub-
446 sequences starting with a known observation, see Methods and Appendix 1.1). This makes
447 the use of MCMC possible in a data-intensive context, in which both the number of free
448 parameters (600 for the biomass matrix, 90 for each of the biodiversity, shade tolerance and
449 stand age matrices) and the number of samples (32,552) constitute increases of several or-
450 ders of magnitude. Similar irregularity problems are quite common in ecological datasets,
451 and the presented approach may have numerous applications beyond the statistical analysis
452 of forest inventories. This methodology can also be applied to other datasets, even with
453 regular samplings, and the same methodology can be applied to deduce transitions with a
454 finer temporal scale.

455 In this study we have analyzed the Quebec forest inventories without explicitly taking
456 into account the geographical location of plots, as well as the environmental and climatic
457 variables. We have obtained transition matrices covering temperate to boreal forests, with
458 a disturbance regime varying from canopy gaps to disastrous fires. We have repeated the

459 developed approach after subdividing the Quebec dataset into the major ecological domains
460 and have not observed substantial differences between the resulting transition matrices and
461 the general matrices presented in this study (Li enard et al. unpublished data). In addition
462 to this, the biomass transition matrices computed for the Lake States in the US (see Strigul
463 et al. (2012) Tables 2 and 3) and the shade tolerance index transition matrices computed in
464 northeastern parts of the US (Lienard et al., 2014) are quite similar to the ones presented
465 in this study. It is quite amazing in fact that we could represent the dynamics of stand
466 level characteristics given the neglect of geography. We hypothesize that the forest stand
467 dynamics as well as disturbance regimes have substantial similarities across a large number
468 of boreal and temperate forest types, and this will be specifically addressed in our future
469 studies. We believe that the ability to make broad predictions on the forest stand dynamics
470 without going into the fine details of geography is one of the major strengths of our approach.

471 The patch-mosaic framework has been already extensively employed in forest model-
472 ing (Kohyama et al., 2001; Kohyama, 2006; Scherstjanoi et al., 2013). Our approach has
473 substantial similarities with the previous studies using the same scientific background (see
474 Appendix 1), however there are distinctions related to the definitions of the forest patch.
475 The forest stand or (forest patch) in this work represent a unit of forest which is large
476 enough to be a community of trees, where individual tree gap dynamics is averaged, but
477 at the same time small enough to be a subject to intermediate and large scale disturbances
478 (Strigul et al., 2012)[p.72]. This definition results in an estimate of about 0.5-1 ha, which
479 allows to use forest inventory permanent plots directly as an approximate forest stand rep-
480 resentations (the size of the standard Quebec forest inventory plot is about 625 m² and the
481 USDA FIA plot is 675 m²). In other application of patch-mosaic concept to forest dynamics
482 the patches (stands) are often defined differently. The size of patches varies from the size
483 of large individual trees (in this case the patch dynamics is essentially equivalent to the gap
484 dynamics Kohyama et al., 2001; Kohyama, 2006; Moorcroft et al., 2001), through patches
485 similar to employed in our study (Acevedo et al., 2001) to the much large forest patches
486 representing many hectares of forest (Boychuk et al., 1997). The difference in definitions of

487 the patch essentially reflects the different applications and questions that can be addressed
488 with particular models (see Strigul et al. (2012)[p.71] for an additional discussion).

489 *5.3. Predictions for forest dynamics in Quebec*

490 Our model made several notable predictions about future forest dynamics in Quebec.
491 The most pronounced predicted changes are substantial short-term increase in biomass and
492 a longer-term increase in average age of trees (Fig. 3). The increase in biomass is intuitively
493 consistent with the increase in stand age, and both demonstrate a progression toward more
494 mature stands. This progression is to be sustained throughout the next 20 years and beyond
495 (Fig. 3), thus meaning that the unmanaged forests sampled in the inventory are currently
496 far from their equilibrium state. The model also predicted smaller changes in biodiversity
497 and the shade tolerance index. To understand stand maturation occurring with the small
498 increase in the prevalence of early successional species, we must recall that neither biomass
499 nor stand age are significantly correlated with shade tolerance index in the dataset (e.g.,
500 $r = -0.02$ with 95% confidence interval [-0.03,-0.01] for biomass and shade tolerance, see
501 Fig. 5 in Appendix). Thus, it is unsurprising that the predictions are not correlated.
502 Further, the predicted changes happen with different temporal dynamics and have different
503 magnitudes, and have probably distinct mechanisms. In particular, while biomass and stand
504 age are affected by both individual tree growth (leading to an increase) and disturbances
505 (leading to a decrease), the shade tolerance index is affected only by disturbances. On
506 the one hand, small disturbances (e.g., individual tree mortality) will typically promote
507 the recruitment of late successional species into the canopy through gap dynamics. On
508 the other hand, intermediate and large-scale disturbances will facilitate early successional
509 species via the development of large canopy openings (e.g. Taylor and Chen, 2011). Thus,
510 increase of intermediate and large-scale disturbances may promote early successional species,
511 while the overall increase in biomass and stand age would result largely from individual tree
512 growth. Our work thus suggests that Quebec forests are not progressing toward higher shade
513 tolerance states despite their continuous biomass and stand age growth. This result echoes
514 recent studies which showed that shade tolerance is not the sole driver for forest succession

515 in Canadian central forests (Taylor and Chen, 2011; Chen and Taylor, 2012).

516 The accurate prediction of the second half of the dataset obtained using only the first half
517 of the dataset demonstrate that the natural disturbance regime in the forest plots sampled
518 in the Quebec inventory did not change substantially over the last 30 years. In the context
519 of global warming, this could mean either that (a) there is no substantial consequence yet
520 on the macroscopic dynamics of Quebec forests or that (b) the climatic change consequences
521 were already present in the first half of the dataset or that (c) our analysis is not fine enough
522 to catch the signal of the recent climate change (in particular moving climatic boundaries,
523 cf. McKenney et al., 2007, are not taken into account as the approach developed here is not
524 spatially explicit). In all cases, the inclusion in the transition matrices of future disturbances
525 induced by climatic change (e.g. the increase of forest fire reviewed in Flannigan et al.,
526 2009) could be a promising follow-up of our work by providing quantitative insights on the
527 consequences of global warming on forests. The study of changes in disturbances was not
528 the focus of the current study, and we have to be careful in generalizing conclusions about
529 global warming as the gradual non-stationary disturbance regimes might take from 50 to
530 100 years to show significant departures (Loudermilk et al., 2013; Rhemtulla et al., 2009;
531 Thompson et al., 2011).

532 The multidimensional nature of forest stands creates substantial challenges for modeling.
533 Our study demonstrates that at least four dimensions are uncorrelated in the Quebec dataset,
534 and that stand characteristics cannot be collapsed around one variable. The data intensive
535 model could be based on uncorrelated principal component axes. However, such a model
536 would not lend itself to a simple mechanistic interpretations in terms of macroscopic forest
537 characteristics. Therefore, we have developed the model using the mostly uncorrelated stand
538 characteristics: biomass, biodiversity, shade tolerance index, and average age of trees. In
539 this model, the small correlations between these characteristics (Figure 3 in Appendix) will
540 propagate to the model predictions, potentially resulting in slightly correlated predictions,
541 in contrast with a model developed on the principal components. However, this choice of
542 dimensional variables has the decisive advantage of allowing for the meaningful interpretation

543 of the transition matrices and predictions.

544 The constructed transition matrices have predictive power, as demonstrated in Section
545 4.3. However, the universality of the predictions is intrinsically dependent on the repre-
546 sentativity of the dataset, and a bias in data collection will be reported in the predictions.
547 For this particular dataset for instance, we observed (and also predicted) increasing biomass,
548 diversity, age and slight decreasing shade tolerance over time. However we should not expect
549 forest stands affected by additional silvicultural operations, such as logging, to follow the
550 trajectory recorded in the Quebec dataset. Thus, predictions made with this dataset should
551 not be extended to them.

552 **Acknowledgements**

553 This work was partially supported by a grant from the Simons Foundation (#283770 to
554 N.S.) and a Washington State University New Faculty SEED grant. D.G. also acknowledges
555 the financial support from a Strategic Grant of NSERC. We thank Matthew Talluto for
556 interesting discussions and help with editing the manuscript.

Acevedo, M., Urban, D., and Shugart, H. (1996). Models of forest dynamics based on roles of tree species. *Ecological Modelling*, 87(13):267 – 284.

Acevedo, M. F., Ablan, M., Urban, D. L., and Pamarti, S. (2001). Estimating parameters of forest patch transition models from gap models. *Environmental Modelling & Software*, 16(7):649 – 658.

Boychuk, D., Perera, A. H., Ter-Mikaelian, M. T., Martell, D. L., and Li, C. (1997). Modelling the effect of spatial scale and correlated fire disturbances on forest age distribution. *Ecological Modelling*, 95(2):145–164.

Caswell, H. (2001). *Matrix population models: construction, analysis, and interpretation*. Sinauer Associates.

Chen, H. Y. and Taylor, A. R. (2012). A test of ecological succession hypotheses using 55-year time-series data for 361 boreal forest stands. *Global Ecology and Biogeography*, 21(4):441–454.

Deltour, I., Richardson, S., and Hesran, J.-Y. L. (1999). Stochastic algorithms for markov models estimation with intermittent missing data. *Biometrics*, 55(2):565–573.

- Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 1–38.
- Easterling, M. R., Ellner, S. P., and Dixon, P. M. (2000). Size-specific sensitivity: applying a new structured population model. *Ecology*, 81(3):694–708.
- Facelli, J. and Pickett, S. T. (1990). Markovian chains and the role of history in succession. *"Trends in Ecology & Evolution "*, 5(1):27 – 30.
- Feller, W. (1971). *An Introduction to Probability Theory and Its Applications*. Wiley.
- Flannigan, M. D., Krawchuk, M. A., de Groot, W. J., Wotton, B. M., and Gowman, L. M. (2009). Implications of changing climate for global wildland fire. *International Journal of Wildland Fire*, 18(5):483–507.
- Gelfand, A. E. and Smith, A. F. (1990). Sampling-based approaches to calculating marginal densities. *Journal of the American statistical association*, 85(410):398–409.
- Geman, S. and Geman, D. (1984). Stochastic relaxation, gibbs distributions, and the bayesian restoration of images. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, PAMI-6(6):721–741.
- Harrison, P. J., Hanski, I., and Ovaskainen, O. (2011). Bayesian state-space modeling of metapopulation dynamics in the glanville fritillary butterfly. *Ecological Monographs*, 81(4):581–598.
- Hill, M. O. (2003). Diversity and evenness: A unifying notation and its consequences. *Ecology*, 54(2):427–432.
- Horn, H. (1981). Some causes of variety in patterns of secondary succession. In West, D., Shugart, H., and Botkin, D., editors, *Forest Succession*, Springer Advanced Texts in Life Sciences, pages 24–35. Springer, New York.
- Horn, H. S. (1974). The ecology of secondary succession. *Annual Review of Ecology and Systematics*, 5:25–37.
- Jenkins, J., Chojnacky, D., Heath, L., and Birdsey, R. (2003). National-scale biomass estimators for united states tree species. *Forest Science*, 49(1):12–35.
- Kelling, S., Hochachka, W., Fink, D., Riedewald, M., Caruana, R., Ballard, G., and Hooker, G. (2009). Data-intensive science: A new paradigm for biodiversity studies. *BioScience*, 59(7):613.
- Kohyama, T. (2006). The effect of patch demography on the community structure of forest trees. *Ecological Research*, 21(3):346–355.
- Kohyama, T., Suzuki, E., Partomihardjo, T., and Yamada, T. (2001). Dynamic steady state of patch-mosaic tree size structure of a mixed dipterocarp forest regulated by local crowding. *Ecological Research*, 16(1):85–98.

- Korotkov, V. N., Logofet, D. O., and Loreau, M. (2001). Succession in mixed boreal forest of russia: Markov models and non-markov effects. *Ecological Modelling*, 142(12):25 – 38.
- Leslie, P. H. (1945). On the use of matrices in certain population mathematics. *Biometrika*, 33:183–212.
- Levin, S. A. (1998). Ecosystems and the biosphere as complex adaptive systems. *Ecosystems*, 1(5):pp. 431–436.
- Levin, S. A. (1999). *Fragile dominion: complexity and the commons*. Perseus Publishing, Cambridge, MA.
- Levin, S. A. (2003). Complex adaptive systems: Exploring the known, the unknown and the unknowable. *American Mathematical Society*, 40:3–19.
- Levin, S. A. and Paine, R. T. (1974). Disturbance, patch formation, and community structure. *Proceedings of the National Academy of Sciences of the United States of America*, 71(7):2744–2747.
- Lienard, J., Florescu, I., and Strigul, N. (2014). An appraisal of the classic forest succession paradigm with the shade tolerance index. <http://dx.doi.org/10.1101/004994>.
- Logofet, D. and Lesnaya, E. (2000). The mathematics of markov models: what markov chains can really predict in forest successions. *Ecological Modelling*, 126(2):285–298.
- Loudermilk, E. L., Scheller, R. M., Weisberg, P. J., Yang, J., Dilts, T. E., Karam, S. L., and Skinner, C. (2013). Carbon dynamics in the future forest: the importance of long-term successional legacy and climate–fire interactions. *Global change biology*, 19(11):3502–3515.
- McKenney, D. W., Pedlar, J. H., Lawrence, K., Campbell, K., and Hutchinson, M. F. (2007). Potential impacts of climate change on the distribution of North American trees. *Bioscience*, 57(11):939–948.
- Michener, W. K. and Jones, M. B. (2012). Ecoinformatics: supporting ecology as a data-intensive science. *Trends in Ecology and Evolution*, 27(2):85 – 93.
- Moorcroft, P., Hurtt, G., and Pacala, S. W. (2001). A method for scaling vegetation dynamics: the ecosystem demography model (ed). *Ecological monographs*, 71(4):557–586.
- Pasanisi, A., Fu, S., and Bousquet, N. (2012). Estimating discrete markov models from various incomplete data schemes. *Computational Statistics & Data Analysis*, 56(9):2609–2625.
- Perron, J., Morin, P., et al. (2011). Normes d’inventaire forestier: Placettes-échantillons permanentes.
- R Core Team (2012). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0.

- Rhemtulla, J. M., Mladenoff, D. J., and Clayton, M. K. (2009). Historical forest baselines reveal potential for continued carbon sequestration. *Proceedings of the National Academy of Sciences*, 106(15):6082–6087.
- Risk, B. B., De Valpine, P., and Beissinger, S. R. (2011). A robust-design formulation of the incidence function model of metapopulation dynamics applied to two species of rails. *Ecology*, 92(2):462–474.
- Robert, C. P. and Casella, G. (2004). *Monte Carlo statistical methods*, volume 319. Citeseer.
- Scherstjanoi, M., Kaplan, J., Thürig, E., and Lischke, H. (2013). Gappard: a computationally efficient method of approximating gap-scale disturbance in vegetation models. *Geoscientific Model Development Discussions*, 6(1):1021–1084.
- Stephens, G. R. and Waggoner, P. E. (1980). A half century of natural transitions in mixed hardwood forests. *Bulletin, Connecticut Agricultural Experiment Station*, 783.
- Strigul, N. and Florescu, I. (2012). Statistical characteristics of forest succession. In *Abstracts of 97th ESA Annual Meeting*. Ecological Society of America.
- Strigul, N., Florescu, I., Welden, A. R., and Michalczewski, F. (2012). Modelling of forest stand dynamics using markov chains. *Environmental Modelling and Software*, 31:64 – 75.
- Strigul, N., Pristinski, D., Purves, D., Dushoff, J., and Pacala, S. (2008). Scaling from trees to forests: Tractable macroscopic equations for forest dynamics. *Ecological Monographs*, 78(4):523–545.
- Taylor, A. R. and Chen, H. Y. (2011). Multiple successional pathways of boreal forest stands in central canada. *Ecography*, 34(2):208–219.
- ter Braak, C. J. and Etienne, R. S. (2003). Improved bayesian analysis of metapopulation data with an application to a tree frog metapopulation. *Ecology*, 84(1):231–241.
- Thompson, J. R., Foster, D. R., Scheller, R., and Kittredge, D. (2011). The influence of land use and climate change on forest biomass and composition in Massachusetts, USA. *Ecological Applications*, 21(7):2425–2444.
- Usher, M. (1969). A matrix model for forest management. *Biometrics*, 25:309–315.
- Usher, M. (1981). Modelling ecological succession, with particular reference to markovian models. In *Vegetation dynamics in grasslands, healthlands and mediterranean ligneous formations*, pages 11–18. Springer, New York.
- Usher, M. B. (1979a). Markovian approaches to ecological succession. *The Journal of Animal Ecology*, pages 413–426.
- Usher, M. B. (1979b). Markovian approaches to ecological succession. *The Journal of Animal Ecology*, 48(2):413–426.

Van Wagner, C. E. (1978). Age-class distribution and the forest fire cycle. *Canadian Journal of Forest Research*, 8(2):220–227.

Waggoner, P. E. and Stephens, G. R. (1970). Transition probabilities for a forest. *Nature*, 225:1160–1161.

Watt, A. S. (1947). Pattern and process in the plant community. *J.Ecol.*, 35:1–22.

Appendix to “Data-intensive multidimensional modeling of forest dynamics”

Jean F. Liénard ¹, Dominique Gravel ², Nikolay S. Strigul ^{1*}

¹: Department of Mathematics, Washington State University, Washington

²: Département de Biologie, Université du Québec à Rimouski, Québec

*: corresponding author email: nick.strigul@vancouver.wsu.edu

Appendix 1.1 Markov chain examples

We illustrate in Fig. 1 the application of the Markov chain methodology to the modeling of forest stand dynamics across three simplified models. All these models have been designed to have the same probability of aging (probability of 20% to go to the next state) and of disturbances (when all earlier states are pooled, the probability to go backward is consistently 10%). However, the long-term distributions are substantially different across models, demonstrating the need for a data-intensive approach that incorporates variable-scale disturbances to quantitatively constrain the transition matrices.

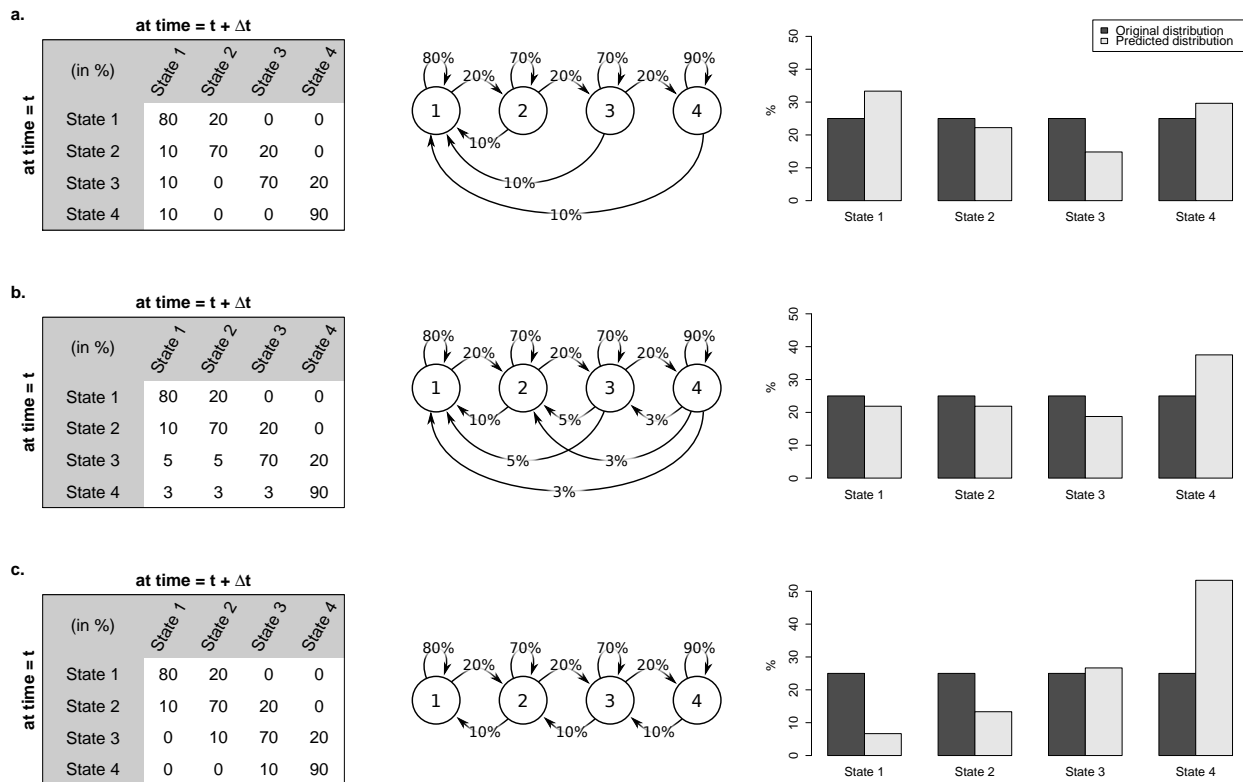


Figure 1: Hypothetical examples illustrating the Markov chain model for forest stand dynamics. Three different scenarios are considered: birth-and-disaster (a), birth-and-disturbances (b) and birth-and-regression (c). For each scenario, we show the transition matrix rounded to the closest percent (left), the corresponding Markov chain graphs (middle) and the long-term equilibrium when starting from the same uniform distribution (right). In all three scenarios, the probability of going to the next state is 20%, and the probability of going to any earlier state is 10%; however, the resulting long-term distribution of states are very different.

Appendix 1.2 Description of the Quebec Dataset

The Quebec forest inventory program was started in 1970s and is still ongoing nowadays, using a constant methodology across the years (Perron et al., 2011). In this study, we considered only permanent plots that were sampled at least twice since the beginning of the inventory, resulting in 32552 plot measurements taken from 11660 different locations throughout Quebec. These permanent plots are circular of area $400m^2$. At the time of establishment as well as at each successive measurements, every tree with a diameter greater than 90 mm was numbered with paint, measured and finally recorded in the database (Perron et al., 2011, Duchesne and Ouimet, 2009). For the computation of the plot characteristics, we specifically relied on the species determined by the forester, the D.B.H. measured using a diameter tape to the nearest millimeter and the age measured from various sources.

The version of the Quebec dataset used in our study spans from 1970 until 2007 included. Fig. 2 shows the distributions of the measurements across the years in the database, as well as the repartition of the intervals between two successive measurements.

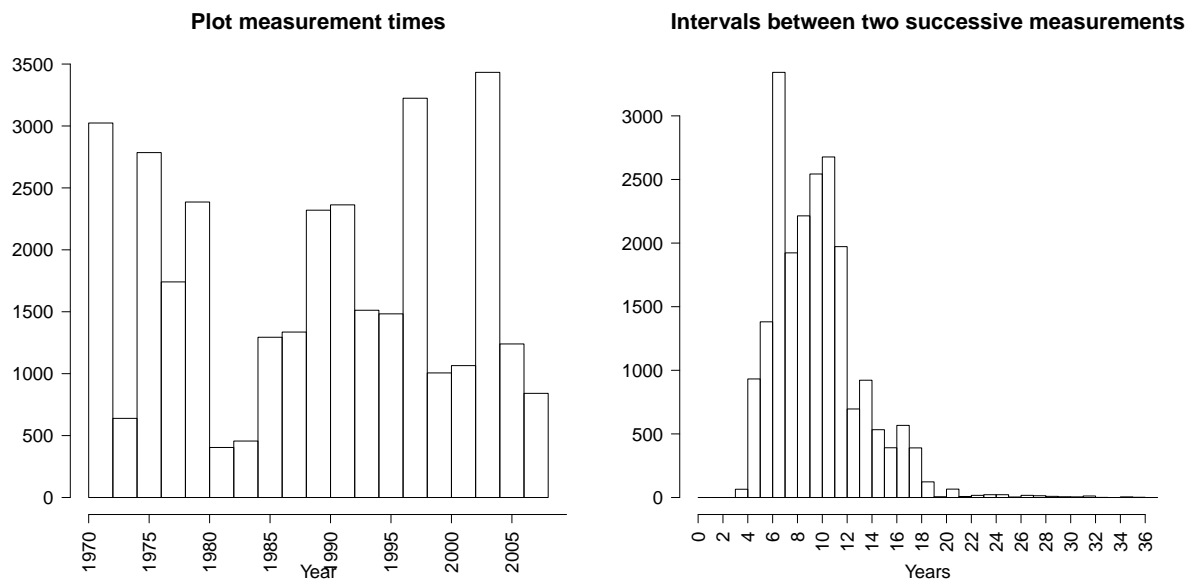


Figure 2: Distribution of the times of plot measurement (left) and distribution of intervals between successive measurements of the same plot (right) in the Quebec dataset.

The distributions of stand characteristics across all the database records is plotted in Fig. 3, and their boxplots are shown in Fig. 4. Different distribution patterns are evident. In particular, the biomass, basal area and stand age distributions reveal the presence of long tails for the high values.

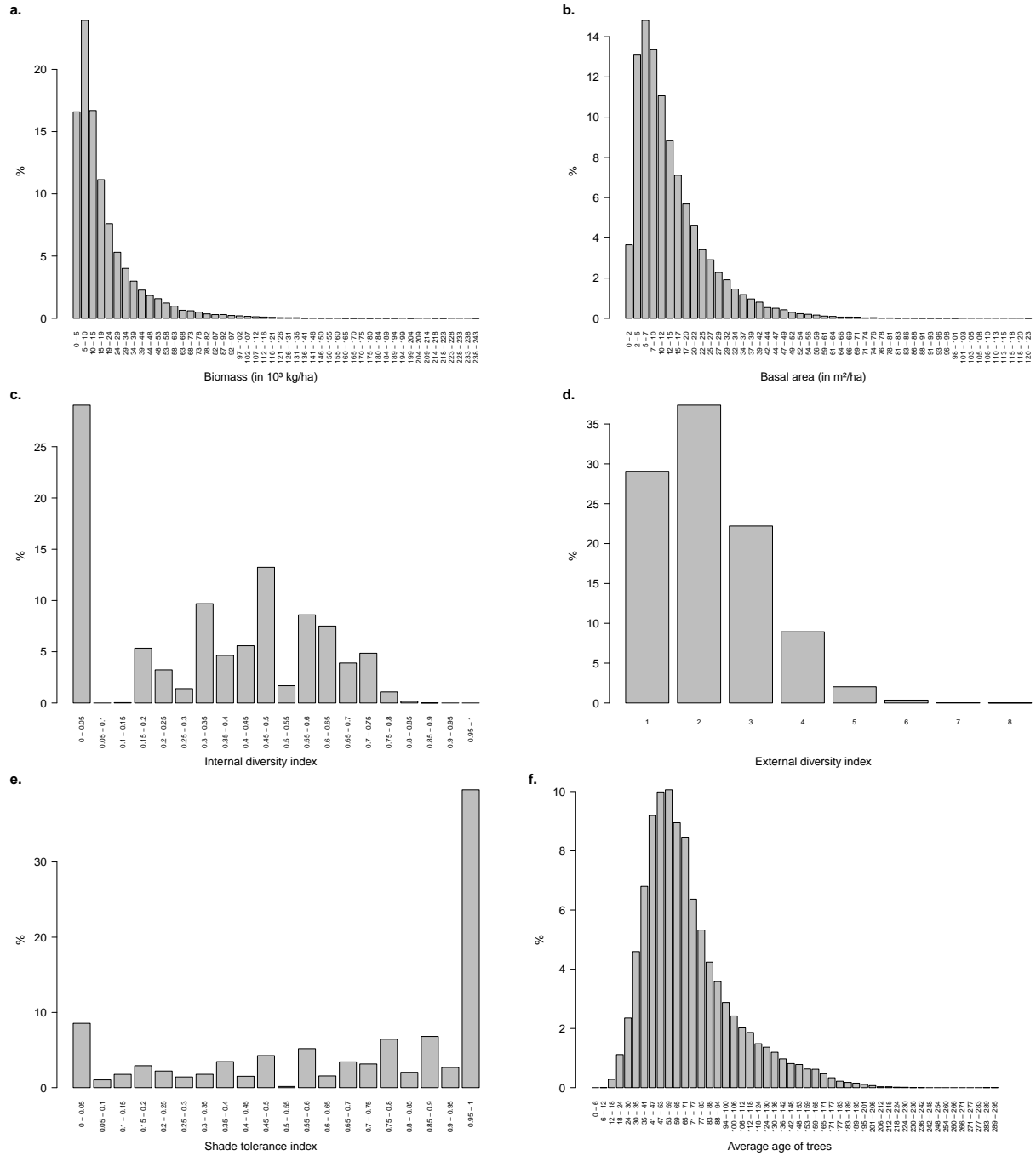


Figure 3: Distribution of stand characteristics

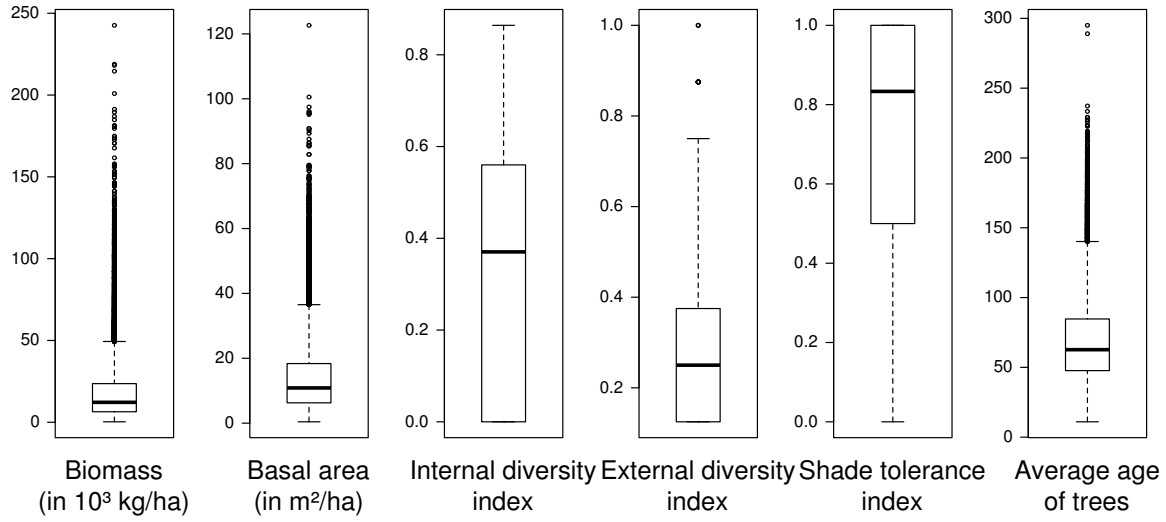


Figure 4: Boxplots of stand characteristics

Appendix 1.3 Gibbs sampling implementation

Forest inventory protocols constitute so-called ignorable data collection mechanisms (Little and Rubin, 1987), as the plots are sampled according to a pre-defined policy that depends mainly on organizational constraints on the foresters' activity, and not on the actual composition of the plots (Perron et al., 2011). This allows for the use of Gibbs sampling as described in main text, which follows the implementation of Pasanisi et al. (2012).

Several steps of data preparation are essential to apply MCMC algorithms to any dataset. We present here two improvements that are specifically relevant for the structure of the Quebec inventory, and which have been used for both the validation and predictions on these data. These are presented and justified briefly for their example values; applying our methodology to different databases could call for alternative reductions, or for no reduction at all if computational resources are not a limiting factor to perform the inference with Gibbs sampling.

First, because no plots were resampled faster than every 3 years in the inventory, we lowered the computational cost by keeping only one state every 3 years. In addition, to further lower the length of sequences, we extracted all sub-sequences of length 4 (hence corresponding to $4 * 3 = 12$ years) containing at least two measurements. These two manipulations shorten the sequences dramatically, from sequences originally spanning over 47 years and composed mostly of missing data, to sequences of length 4 and containing at least 2 measurements.

See the following box for a concrete example encompassing all the steps used to include yearly characteristics into temporal sequences.

We provide as an example the data preparation with biomass modified from the Quebec inventory for one plot:

- 1975: biomass = $8.2 \cdot 10^3$ kg/ha
- 1981: biomass = 4.6
- 1988: biomass = 3.9
- 1992: biomass = 13.7
- 2000: biomass = 14.3

The preliminary step is to express these numerical values into discrete states. In our study, we subdivided biomass into 25 states from 0 to $50 \cdot 10^3$ kg/ha, corresponding to the following description of the same plot's biomass:

- 1975: biomass $\in 8 - 10 \cdot 10^3$ kg/ha
- 1981: biomass $\in 4 - 6$
- 1988: biomass $\in 2 - 4$
- 1992: biomass $\in 12 - 14$
- 2000: biomass $\in 14 - 16$

The first proposed step to simplify this temporal sequence is discretization into 3-year intervals:

70-72	73-75	76-78	79-81	82-84	85-87	88-90	91-93	94-96	97-99	00-02	03-05	06-08
?	8-10	?	4-6	?	?	2-4	12-14	?	?	14-16	?	?

The second step of discretization is to further split the sequence into sub-sequences starting with a known observation and containing at least one additional observation, as well as dismissing the absolute value of the year in order to retain only the relative dynamics:

	T	T+3 years	T+6 years	T+9 years
Sequence 1	8-10	?	4-6	?
Sequence 2	4-6	?	?	2-4
Sequence 3	2-4	12-14	?	?
Sequence 4	12-14	?	?	14-16

The sparse sequences obtained as explained above are used to infer the underlying transition matrix. The precise implementation of Gibbs sampling is summarized in main text and summarized there with a box containing the pseudo-code. Informally, the core idea of our implementation of the Gibbs sampling algorithm is to estimate both the transition matrix and the missing values of these sequences. To do this, we perform two iterative steps, which we refer to as “parameter estimation” and “data augmentation”. In the parameter

estimation step, we sample transition probabilities according to the current estimate of the missing data. This results in an estimate of the transition matrix. In the data augmentation step, we sample values for filling the missing measurements based on the current estimate of the transition matrix. This results in an estimate of the missing data. As these two steps are interdependent, starting the algorithm is a “chicken and egg” type of problem, which is solved with a random initialization of the missing states and the subsequent dismissal of the first B iterations (Robert and Casella, 2004, Pasanisi et al., 2012).

Appendix 1.4 Correlation Analysis

Correlation matrices using the Pearson correlation coefficient were computed for all years (Fig. 5), as well as decade-by-decade (see the following tables):

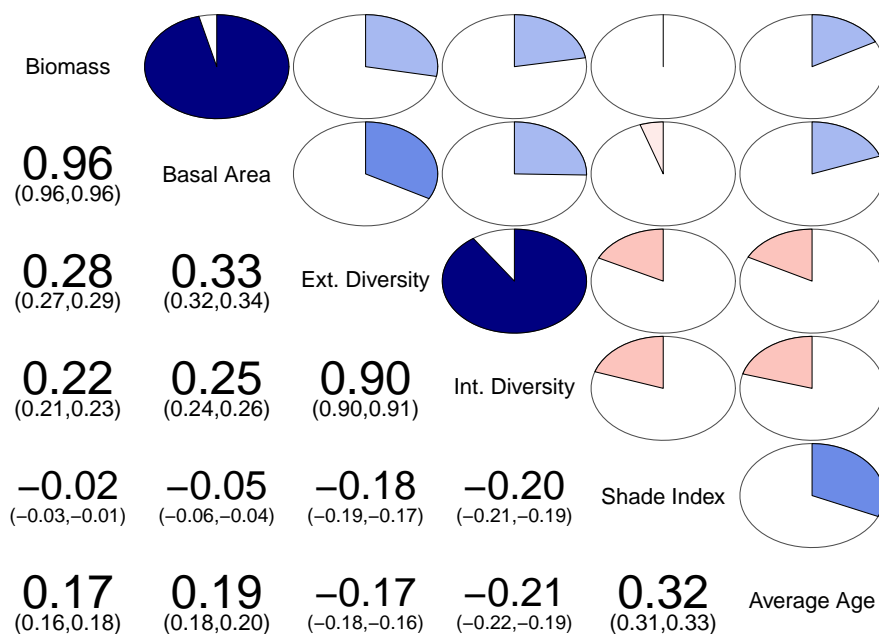


Figure 5: Graphical display of the correlation matrix for each characteristic; the numbers below the diagonal are the Pearson's coefficients (with 95% confidence intervals) and the circles above the diagonal provide visual indications of the correlations, with blue colors denoting positive correlations and pink colors denoting negative correlations (Friendly, 2002).

Correlation for 2000s

age-supplied dataset (6292 data)	BIOMASS	BA	MEANSUCC	EXTDIV	INTDIV	MEANAGE
BIOMASS	1.00	0.96	-0.11	0.36	0.33	0.08
BA	0.96	1.00	-0.15	0.40	0.37	0.08
MEANSUCC	-0.11	-0.15	1.00	-0.24	-0.27	0.35
EXTDIV	0.36	0.40	-0.24	1.00	0.89	-0.15
INTDIV	0.33	0.37	-0.27	0.89	1.00	-0.16
MEANAGE	0.08	0.08	0.35	-0.15	-0.16	1.00

Correlation for 1990s

age-supplied dataset (9845 data)	BIOMASS	BA	MEANSUCC	EXTDIV	INTDIV	MEANAGE
BIOMASS	1.00	0.96	-0.02	0.31	0.28	0.10
BA	0.96	1.00	-0.05	0.37	0.33	0.10
MEANSUCC	-0.02	-0.05	1.00	-0.15	-0.18	0.35
EXTDIV	0.31	0.37	-0.15	1.00	0.90	-0.22
INTDIV	0.28	0.33	-0.18	0.90	1.00	-0.24
MEANAGE	0.10	0.10	0.35	-0.22	-0.24	1.00

Correlation for 1980s

age-supplied dataset (6097 data)	BIOMASS	BA	MEANSUCC	EXTDIV	INTDIV	MEANAGE
BIOMASS	1.00	0.96	0.06	0.18	0.14	0.11
BA	0.96	1.00	0.06	0.22	0.17	0.17
MEANSUCC	0.06	0.06	1.00	-0.16	-0.17	0.40
EXTDIV	0.18	0.22	-0.16	1.00	0.92	-0.20
INTDIV	0.14	0.17	-0.17	0.92	1.00	-0.21
MEANAGE	0.11	0.17	0.40	-0.20	-0.21	1.00

Correlation for 1970s

age-supplied dataset (9417 data)	BIOMASS	BA	MEANSUCC	EXTDIV	INTDIV	MEANAGE
BIOMASS	1.00	0.95	-0.02	0.22	0.18	0.26
BA	0.95	1.00	-0.03	0.26	0.19	0.29
MEANSUCC	-0.02	-0.03	1.00	-0.20	-0.23	0.26
EXTDIV	0.22	0.26	-0.20	1.00	0.92	-0.16
INTDIV	0.18	0.19	-0.23	0.92	1.00	-0.19
MEANAGE	0.26	0.29	0.26	-0.16	-0.19	1.00

Correlation for all years

age-supplied dataset (31651 data)	BIOMASS	BA	MEANSUCC	EXTDIV	INTDIV	MEANAGE
BIOMASS	1.00	0.96	-0.04	0.29	0.23	0.17
BA	0.96	1.00	-0.07	0.33	0.25	0.19
MEANSUCC	-0.04	-0.07	1.00	-0.18	-0.21	0.32
EXTDIV	0.29	0.33	-0.18	1.00	0.90	-0.17
INTDIV	0.23	0.25	-0.21	0.90	1.00	-0.21
MEANAGE	0.17	0.19	0.32	-0.17	-0.21	1.00

Appendix 1.5 Principal component analysis

To confirm the correlation analysis results, we applied a principal component analysis to the dataset, summarized here:

principal component analysis	PC1	PC2	PC3	PC4	PC5	PC6
Standard deviation	1.589	1.3124	0.9958	0.7927	0.30958	0.19001
Proportion of Variance	0.421	0.2871	0.1653	0.1047	0.01597	0.00602
Cumulative Proportion	0.421	0.7080	0.8733	0.9780	0.99398	1.00000

The loadings are presented in the following table and can be visualized in the biplots (Fig. 6).

Loadings	PC1	PC2	PC3	PC4	PC5	PC6
Bionass	0.479	0.43	-0.24	0.200	-0.142	0.683
Basal area	0.495	0.42	-0.23	0.131	0.091	-0.710
Shade tolerance index	-0.171	0.35	0.70	0.593	-0.011	-0.029
External diversity	0.509	-0.31	0.36	-0.082	0.700	0.133
Internal diversity	0.484	-0.36	0.37	-0.102	-0.694	-0.099
Average tree age	-0.042	0.54	0.36	-0.758	-0.012	0.035

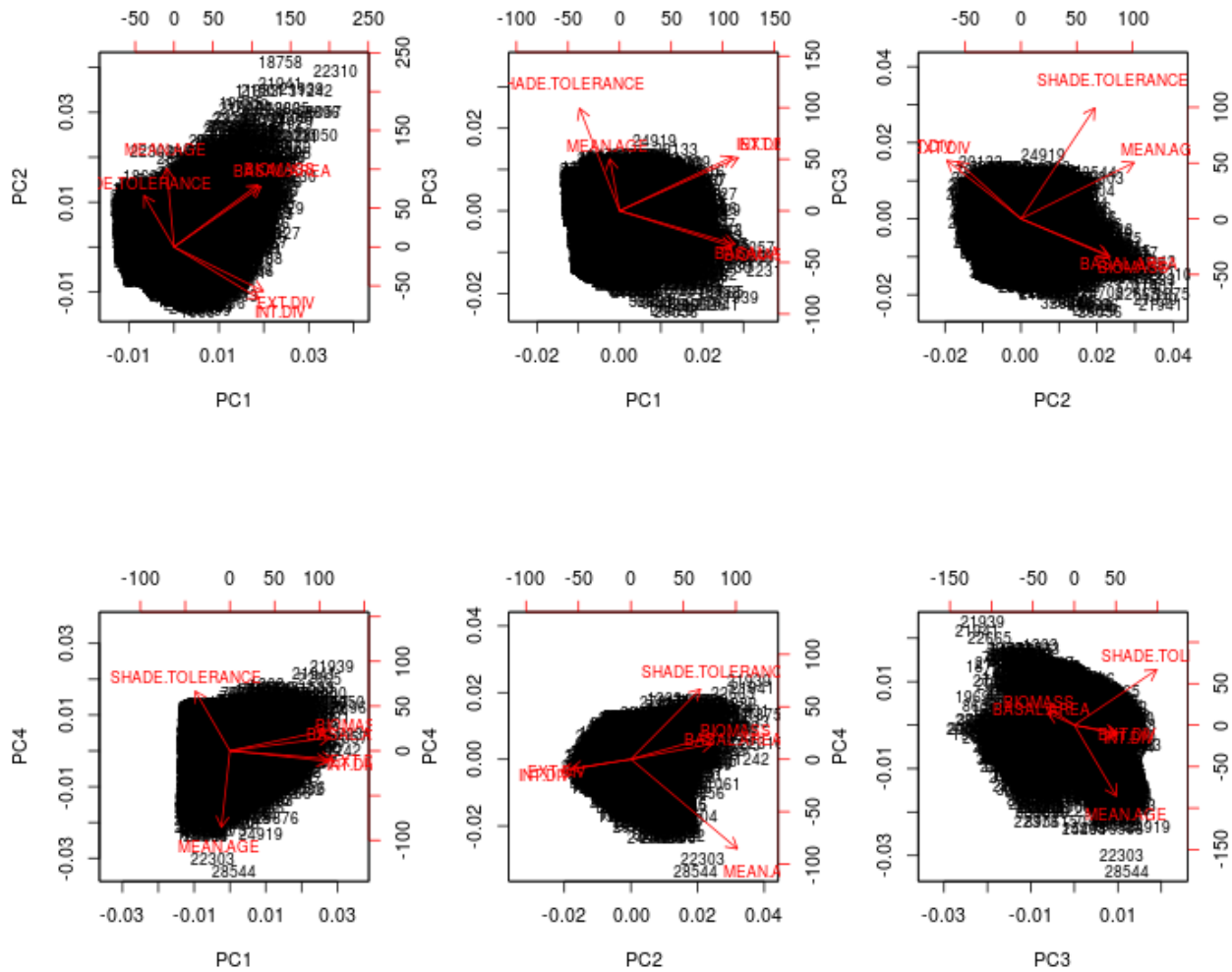


Figure 6: Biplots of all six studied characteristics in the spaces defined by the first four principal components (PC1 to PC4).

Appendix 1.6 Transition matrices of biodiversity, shade tolerance index and average age

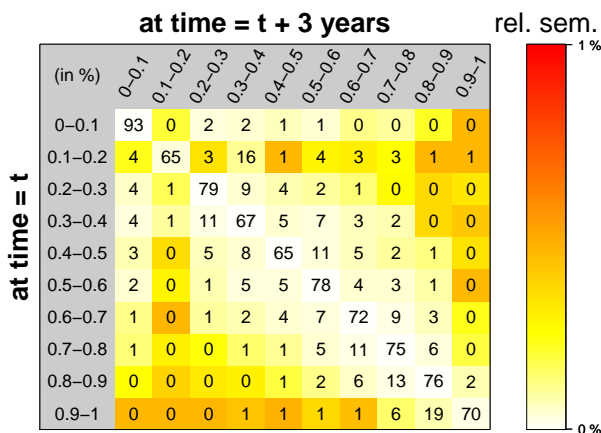


Figure 7: 3-year transition matrix of the internal diversity.

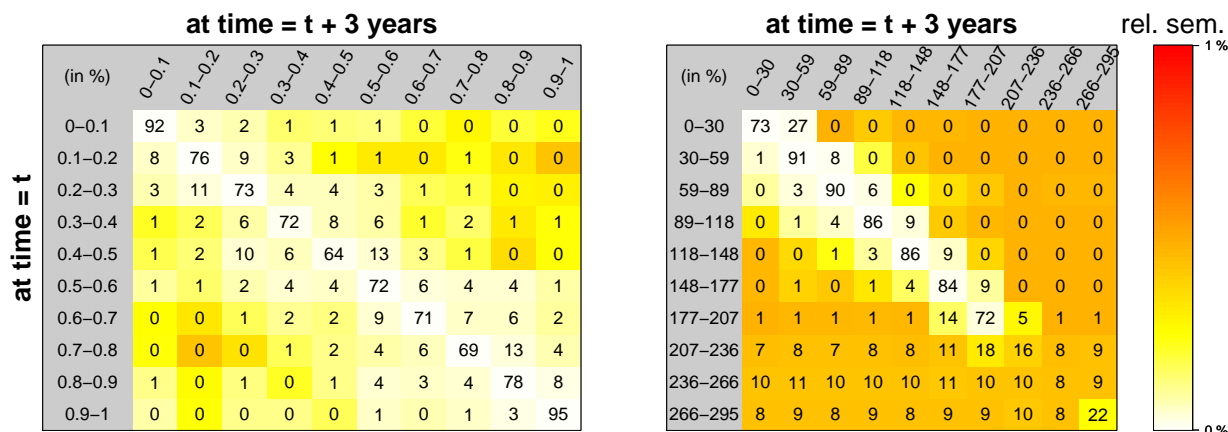


Figure 8: 3-year transition matrices of the shade tolerance index (left) and stand age (right).

Appendix 1.7 Validation

To estimate the general errors in the ability of a Markov-Chain model to predict future forest characteristics based on a dataset such as the Quebec dataset, we performed two different validations by splitting the dataset using two different partitions.

In the first validation, we ran the Gibbs sampler with only the first 18 years of records from 1970 to 1988. We then used the obtained models to predict what would be the forest state for a period corresponding to the second half of the dataset. In details, we pooled all data from the early years and computed from them an aggregated distribution of the first half of the dataset. We then predicted the aggregated distribution 18 years later by iterating 6 times in our 3-year transition models. We finally compared this predicted average with the aggregated distribution in the dataset of the years spanning from 1989 to 2007 (Fig. 3 in main text). As mentioned in the article, the predictions were highly accurate, with R^2 between observation and prediction ranging from 0.8 to 0.95.

In the second validation, we randomly split the data in two different halves, regardless of the year. We then proceeded to the computation of the transition matrix for each half, and computed the corresponding equilibrium states (Fig. 9).

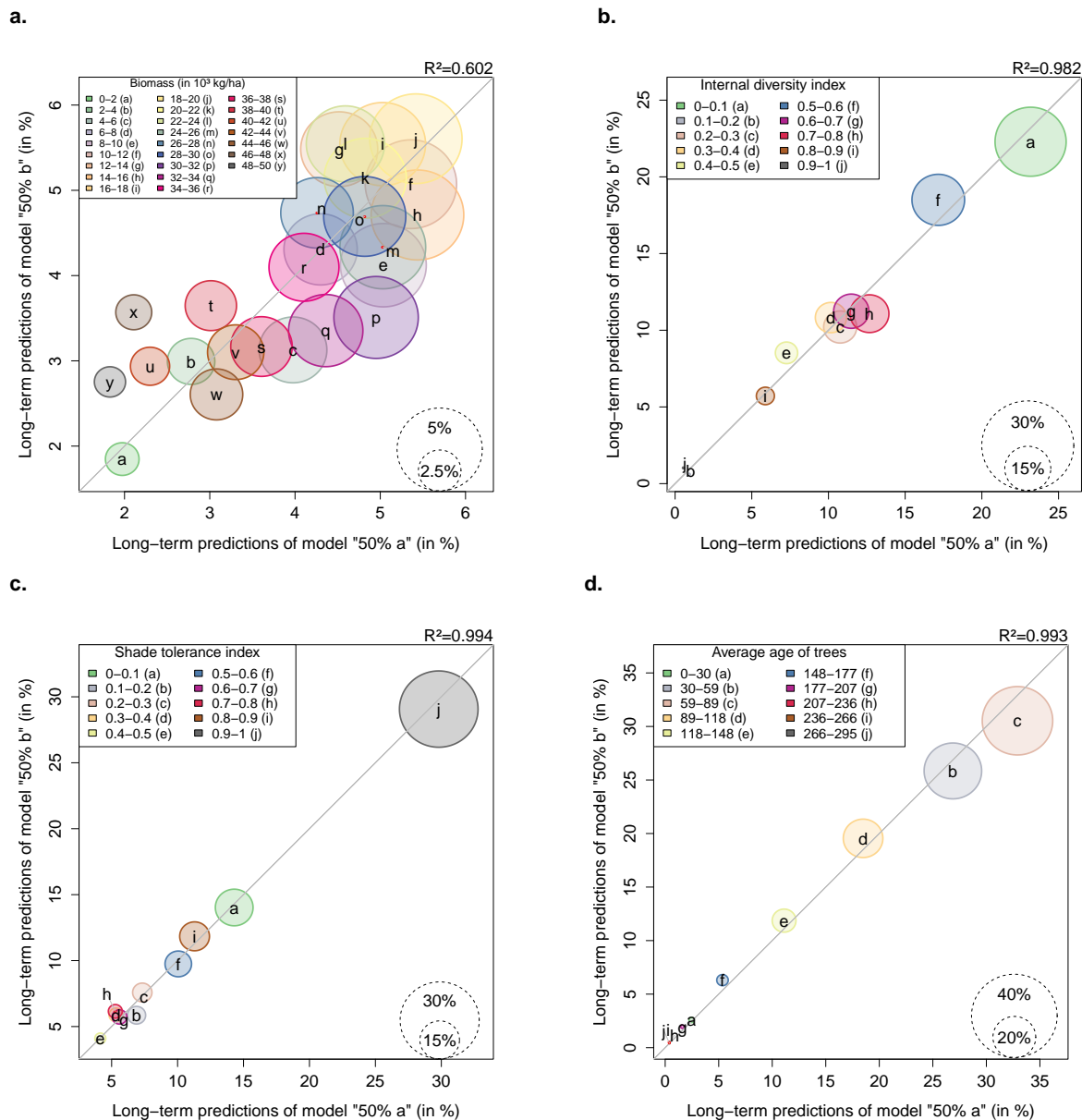


Figure 9: Predictions at equilibrium from two models, each computed with a different half of the dataset (denoted here “50% a” and “50% b”) for the states of each characteristic. For each state, the circle size denotes the number of stands belonging to it in the real dataset. The R^2 error measure is indicated on the top right of each plot. In this validation, the two halves of data used to compute the transition matrix correspond to a random split of the dataset, regardless of the years (see text for details).

References

- Duchesne, L. and Ouimet, R. (2009). Relationships between structure, composition, and dynamics of the pristine northern boreal forest and air temperature, precipitation, and soil texture in quebec (canada). *International Journal of Forestry Research*, 2009.
- Friendly, M. (2002). Corrgrams: Exploratory displays for correlation matrices. *The American Statistician*, 56(4):316–324.
- Little, R. J. and Rubin, D. B. (1987). *Statistical analysis with missing data*, volume 539. Wiley New York.
- Pasanisi, A., Fu, S., and Bousquet, N. (2012). Estimating discrete markov models from various incomplete data schemes. *Computational Statistics & Data Analysis*, 56(9):2609–2625.
- Perron, J., Morin, P., et al. (2011). Normes d’inventaire forestier: Placettes-échantillons permanentes.
- Robert, C. P. and Casella, G. (2004). *Monte Carlo statistical methods*, volume 319. Citeseer.