# DETECTING AND MODELING CHANGES IN A TIME SERIES OF PROPORTIONS

BY THOMAS J. FISHER[1,a], JING ZHANG[1,b], STEPHEN P. COLEGATE[2,c] AND
MICHAEL J. VANNI[3,d]

[1]*Department of Statistics, Miami University,* [a]*fishert4@miamioh.edu,* [b]*zhangj8@miamioh.edu*

[2]*Division of Biostatistics & Bioinformatics, Department of Environmental & Public Health Sciences, University of Cincinnati,*
[c]*colegasn@mail.uc.edu*

[3]*Department of Biology, Miami University,* [d]*vannimj@miamioh.edu*

We propose a framework to detect and model shifts in a time series of continuous proportions, that is, a vector of proportions measuring the parts of a whole. By reparameterizing the shape of a Dirichlet distribution, we can model the location and scale separately through generalized linear models. A hidden Markov model allows the coefficients of the generalized linear models to change, thus allowing for the time series to undergo multiple regimes. This framework allows a practitioner to adequately model seasonality, trends, or include covariate information as well as detect change points. The model's behavior is studied via simulation and through the analysis of lake phytoplankton data from 1992 through 2012. Our analyses demonstrate that the model can be effective in detecting and modeling changes in a time series of proportions. Pertaining to the phytoplankton data, the overall biomass has grown with some changes to the community level dynamics occurring circa 2000. Specifically, the proportion of cyanobacteria appears to have increased to the detriment of diatoms.

**1. Introduction.** Phytoplankton are microscopic, autotrophic organisms found in oceans, seas, and freshwater basin ecosystems. They typically live near surface waters where light is sufficient for growth, as single cells or as colonies that can be visible to the naked eye. A key component of the food web, they transform energy via photosynthesis from sunlight to organic matter that provides food for other organisms. Not all phytoplankton are to be treated equal, however, especially in an environmental context. In 2014, a bloom of toxic cyanobacteria (or blue-green algae) contaminated Lake Erie near Toledo, Ohio, USA, shutting down the city's supply of drinking water. Likewise, a bloom of diatoms in the Pacific northwest affected fisheries in 2015; an increase in domoic acid in shellfish and other small marine animals can be toxic for larger invertebrates. Harmful algal blooms are increasing in frequency all over the world (Paerl, Otten and Kudela (2018)).

Acton Lake is a eutrophic reservoir with high concentrations of nutrients, inorganic sediments, and phytoplankton located in Hueston Woods State Park in southwestern Ohio. Since the early 1990s, the levels of sediment and nutrients entering the lake have fluctuated (Renwick et al. (2018)), along with a growth in phytoplankton abundance (Kelly et al. (2018)). Since 1994, water samples have regularly been recorded for various measurements of water quality, species abundance, and environmental aspects (Kelly et al. (2018)). Water samples are collected regularly from late April through October while access to the study site is available (the marina is closed during the winter).

We study the abundance of phytoplankton taxa, along with environmental conditions (e.g., water temperature), collected over a 21-year period with approximately 11 to 13 measurements per year during the nonwinter months. Data are collected in two-week intervals with
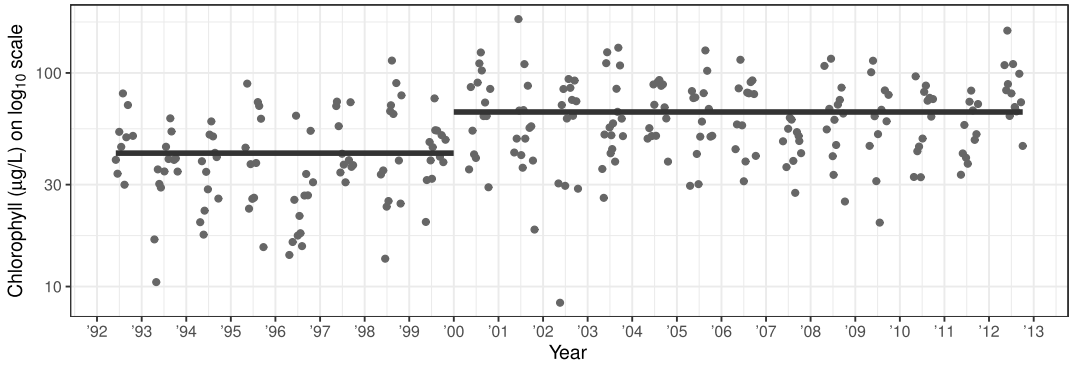
FIG. 1.    *Raw measures of chlorophyll, on log-10 scale, with mean-shift segmentation (based on aggregated data) suggested by change point test in Robbins et al. (2011a).*

the occasional missing value or longer collection interval. Given the irregular time series, we aggregate our measurements into three measures per year, corresponding to the late-spring mixing period, summer stratification, and fall mixing periods: measurements before June 16 are considered *Spring*, from June 16 through July 31 is *Summer*, and measurements from August 1 or later are *Fall*. Other levels of aggregation were considered, but the results are similar (see Supplementary Material, Fisher et al. (2022)), so we only report the results based on three measurements per year here.

Phytoplankton biomass was estimated using the concentration of chlorophyll, as this is the most direct measure of phytoplankton abundance (Kelly et al. (2018)). Figure 1 displays the measurements of chlorophyll at our study site on a log-10 scale along with the segmentation suggested by the ARMA residual change point test from Robbins et al. (2011a) on the aggregate data (deseasoned geometric mean of $\log_{10}$(chlorophyll) with an AR(1) correlation structure): test stat: 1.6507, change point time: fall 1999, *p*-value: 0.0086.

Figure 1 shows an increase in phytoplankton abundance but provides no information on the dynamics of different phytoplankton groups (taxa). For taxonomic information the biovolume of individual species (or subsets of species) was collected by manually identifying, counting, and measuring cells in water samples using a microscope (Hayes and Vanni (2018)). This task is arduous and somewhat susceptible to variability among human counters, particularly for species-level identifications. However, identification at the level of taxonomic groups is more feasible and less subject to human error. Phytoplankton were identified to the finest possible taxonomic resolution but are aggregated into four taxonomic groups: diatoms, flagellates (mostly cryptomonads but also dinoflagellates and euglenoids), chlorophytes ("green algae") and cyanobacteria ("blue-green algae"). Figure 2 displays the phytoplankton compositions in time faceted by season. Of particular interest is how the proportion of phytoplankton groups are changing in time and what (if any) external variables may be driving that change.

1.1. *Compositional data.*    Following its primitive identification in Pearson (1897), the analysis of compositional data, a multivariate set constrained such that each element is non-negative and the set of variables sums to one, owes much of its modern development to Aitchison (1982, 1985). A random observation (or composition) of dimension *p*, $\mathbf{Y}_t$, can be expressed as

$$\mathbf{Y}_t = (Y_1, Y_2, \ldots, Y_p)'_t,$$

$$\text{subject to}\quad Y_i \geq 0, \quad i = 1, \ldots, p, \quad \text{and}\quad \sum_{i=1}^{p} Y_i = 1,$$
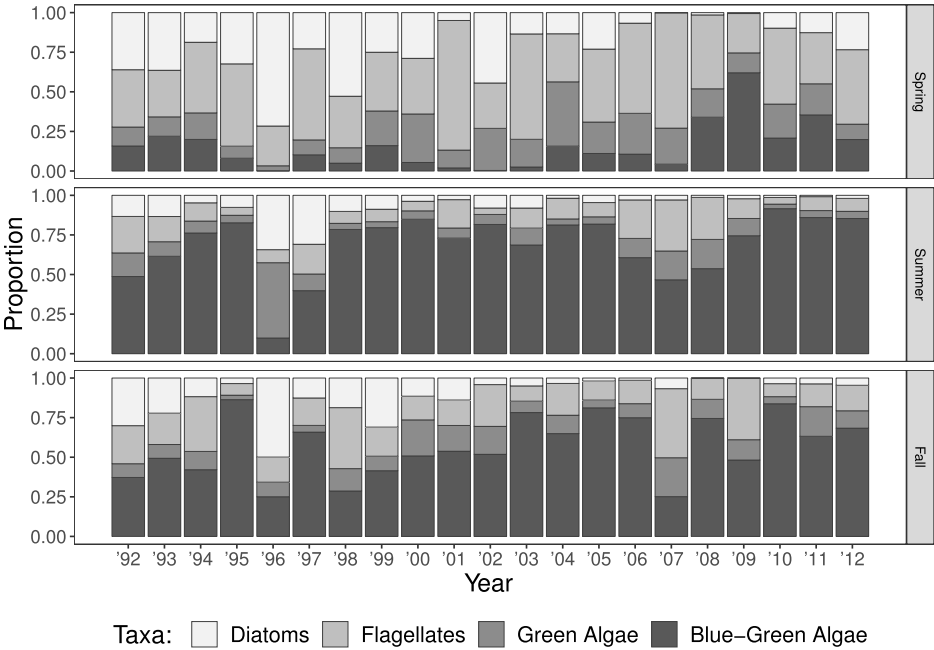
FIG. 2. *Composition of phytoplankton in time, faceted by seasons. There are three measurements per year over 21-years of four-dimensional data. Note, the time series has seasonal effects, and there appears to be a change in the proportions during the length of study.*

where $\mathbf{A}'$ is the transpose of vector/matrix $\mathbf{A}$. It is typically described as a set of measurements representing the parts of a whole, disclosing relative information. In the present application we consider the composition of phytoplankton abundance separated into taxomic groups. Mathematically, observations such as $\mathbf{Y}_t$ fall inside the simplex of dimension $p$, defined as

$$\mathcal{S}^p = \left\{ \mathbf{y} = (y_1, y_2, \ldots, y_p) : y_i > 0 (i = 1, 2, \ldots, p), \sum_{i=1}^{p} y_i = 1 \right\}.$$

When $p = 2$, we are essentially working with a univariate proportion since $Y_2 = 1 - Y_1$. When $p = 3$, the simplex can be visualized through a ternary plot, a particular case of the more general barycentric plot. At higher dimensions, visualizations of this constrained subspace of $\mathbb{R}^p$ is difficult.

Since the seminal work of Aitchison (1986), the development of methods for compositional data has continued to grow but has been sporadic compared to other areas. In part, this is because a key method in compositional data analysis is to project a $p$-dimensional composition into $p - 1$ space via a log-ratio transformation. From there, standard multivariate methods are utilized on the transformed data (see Pawlowsky-Glahn and Buccianti (2011)).

Work on multivariate time series has continued to expand (see Binder, Pourahmadi and Mjelde (2018), Matteson and Tsay (2011)) and is known to have applications in the field of ecology (e.g., Hampton et al. (2013)) and economics (Tsay (2010)). Historically, the log-ratio approach would be applied to time series compositions and then vector autoregressive or state-space models would be utilized (Barceló-Vidal, Aguilar and Martín-Fernández (2011)). Grunwald, Raftery and Guttorp (1993) use a state space approach, and, through reparameterizing of the Dirichlet distribution, they provide a framework that allows the incorporation of predictor variables.

Often Bayesian methods are utilized in multivariate time series, due to estimation inefficiencies in higher dimensions, and to induce necessary structure (Koop and Korobilis (2010),

West (2020)). In fact, in Grunwald, Raftery and Guttorp (1993), a Bayesian approach is used for the model development, although they use a maximum likelihood technique for model estimation.

1.2. *Change points*.    The area of change-point analysis continues to receive interest, particularly in time series (see Aue and Horváth (2013), Bardwell and Fearnhead (2017), Robbins et al. (2011b), e.g.). There, an abrupt change to the distribution of observations is detected at some (unknown) time point, segmenting the time series into regimes. This may be as simple as a shift in mean but may also be in terms of coefficients in a regression model (Lund and Reeves (2002), Robbins, Gallagher and Lund (2016)). Modern methods look to detect change points in multivariate data (Holmes, Kojadinovic and Quessy (2013)). Matteson and James (2014b) propose a distribution free method for detecting multiple change points in multivariate series that has good asymptotic properties. Recently, Prabuchandran et al. (2021) proposed a test for compositional data, based on permutation methods of the log-likelihood value under an i.i.d. Dirichlet assumption.

In a Bayesian framework, much work has been done on the change-point problem as well. Carlin, Gelfand and Smith (1992) considered a hierarchical approach to find a single change point and derive the conditionals of the posterior distribution under certain distributional assumptions. Stephens (1994) extended the Gibbs sampler in Carlin, Gelfand and Smith (1992) to the case of multiple change points. In Barry and Hartigan (1993) the observations are segmented by a product partition model and that approach can also find multiple change points. Erdman and Emerson (2008) improve its implementation and apply it to a large multivariate microarray data set. Kang et al. (2018) develop a Bayesian approach to look for changes in the variance of a regression model, and Liang et al. (2019) use Bayesian change-point methods to explore the relationship of phytoplankton and nutrient levels. The Bayesian framework provides flexible model structures to detect multiple changes in terms of model parameters. However, the implementation of Bayesian methods remain less accessible to data practitioners, especially when the responses are not from an underlying Gaussian process.

The hidden Markov model (HMM) is a valuable tool in the statistical arsenal and has been shown to have many applications. HMMs have a history in change-point detection and modeling (see Chib (1998), Fearnhead (2006), Fearnhead and Liu (2007), Luong, Perduca and Nuel (2012), for examples). In HMMs the distribution of an observation at time $t$ depends on an unobserved (latent) state. The latent states are assumed to follow a Markov process, and, in general, the underlying Markov process allows the distribution of observations to transition between states. When used for change-point detection, a constrained HMM may be implemented, as the underlying data is not necessarily assumed to move back-and-forth between states but rather transition into unique regimes.

1.3. *Contribution*.    This article considers detecting and modeling changes in a time series of continuous proportions and modeling those changes with covariate information. Although each of those topics have been studied marginally, connecting all the aforementioned areas appears to be an understudied problem. Several of the methods in the literature can detect multiple change points on multivariate data but are not necessarily designed for data in the simplex. Prabuchandran et al. (2021) propose a change-point test for compositional data, but it is not designed to work with seasonal data or potential predictor variables. These approaches may find a change point, but they do not simultaneously model the process at the same time.

The outline of this article is as follows: Section 2 describes our proposed methodology where a HMM with a Dirichlet regression detects and models changes in a time series of proportions. The HMM not only can detect changes but also quantify the probability of a

change at a given time. The Dirichlet regression allows a practitioner to separately model the underlying mean and variance of the composition. Section 3 provides a simulation study, demonstrating the efficacy of our method. An in-depth analysis of the phytoplankton data and some conclusions are provided in Sections 4 and 5.

**2. Methodology.** The foundation is the Dirichlet distribution which is considered the most natural distribution for working with compositions (Grunwald, Raftery and Guttorp (1993)). Unlike the log-ratio approaches in Aitchison (1986) and used in other time series applications (see Brunsdon and Smith (1998), Mills (2010), to name a few), working directly in the simplex allows for an easy interpretation of the behavior of the components. For random composition $\mathbf{Y}_t$, the Dirichlet distribution is defined as

$$f(\mathbf{y}|\boldsymbol{\alpha}) = \mathcal{D}(\boldsymbol{\alpha}) \prod_{i=1}^{p} y_i^{\alpha_i - 1},$$

where $\mathcal{D}(\boldsymbol{\alpha})$ is the Dirichlet function on the vector of shape parameters $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \ldots, \alpha_p)$, $\alpha_i \in \mathrm{R}^+$, for $i = 1, 2, \ldots, p$, defined as

$$\mathcal{D}(\boldsymbol{\alpha}) = \frac{\prod_{i=1}^{p} \Gamma(\alpha_i)}{\Gamma(\alpha_1 + \alpha_2 + \cdots + \alpha_p)}$$

and $\Gamma(\cdot)$ is the Gamma function. The domain of the Dirichlet distribution is the simplex of dimension $p$, $\mathcal{S}^p$, and the Dirichlet distribution is a multivariate generalization of the common Beta distribution. It is well known that the expectation and variance of $Y_i$, the $i$th component of $\mathbf{Y}$, is

$$E[Y_i|\boldsymbol{\alpha}] = \alpha_i/\boldsymbol{\alpha}'\mathbf{1}_p \quad \text{and} \quad \mathrm{Var}[Y_i|\boldsymbol{\alpha}] = \frac{\alpha_i(\boldsymbol{\alpha}'\mathbf{1}_p - \alpha_i)}{(\boldsymbol{\alpha}'\mathbf{1}_p)^2(\boldsymbol{\alpha}'\mathbf{1}_p + 1)},$$

where $\mathbf{1}_p$ is a $p$-dimensional vector of ones. The shape parameter $\alpha_i$ essentially controls the behavior of component $Y_i$ in $\mathbf{Y}$.

2.1. *Generalization of Dirichlet assumption.* The shape parameter of the Dirichlet distribution, $\boldsymbol{\alpha}$, controls both the location and dispersion of the distribution. If $\alpha_i$ were to increase while the other terms are held constant, we would expect the proportion of component $i$ to increase and the dispersion of that component to decrease. To isolate the effects of the location and dispersion, we use the approach of Grunwald, Raftery and Guttorp (1993) by reparameterizing the distribution with location, $\boldsymbol{\theta}$, and scale, $\tau$, parameters. This approach relaxes some of the independence properties associated with the Dirichlet distribution; the term $\tau$ not only influences the variance of component $i$ but also the covariance of components $Y_i$ and $Y_j$, $i \neq j$. Let $\boldsymbol{\theta} = \boldsymbol{\alpha}/\tau$, where $\tau = \boldsymbol{\alpha}'\mathbf{1}_p$, thus $\mathbf{Y} \sim \mathrm{Dirichlet}(\boldsymbol{\alpha} = \tau\boldsymbol{\theta})$, with

$$E[\mathbf{Y}|\boldsymbol{\theta}, \tau] = \boldsymbol{\theta} \quad \text{and} \quad \mathrm{Var}[\mathbf{Y}|\boldsymbol{\theta}, \tau] = \boldsymbol{\theta}\boldsymbol{\theta}'/(\tau + 1).$$

The parameter space for the location parameter $\boldsymbol{\theta}$ is the simplex of dimension $p$, $\mathcal{S}^p$, and determines the mean of $\mathbf{Y}$. The scale parameter $\tau$ is a strictly positive value, has no influence on the expectation, and only influences the underlying variability (inversely) and correlation between components. This reparameterization allows us to separately model the location and scale.

The parameter $\boldsymbol{\theta}$ is in $\mathcal{S}^p$, and any estimation of it must account for the constrained nature of the simplex. As in Grunwald, Raftery and Guttorp (1993), we use a Bayesian approach for model development and induce a Dirichlet prior distribution on $\boldsymbol{\theta}$; that is, assume $\boldsymbol{\theta} \sim \mathrm{Dirichlet}(\boldsymbol{\eta})$. The expectation of $\boldsymbol{\theta}$ is $E[\boldsymbol{\theta}] = \boldsymbol{\eta}/(\boldsymbol{\eta}'\mathbf{1}_p)$, and we can model the location using a generalized linear model on the $\boldsymbol{\eta}$ parameter (Hijazi and Jernigan (2009)).

The scale parameter $\tau$ is strictly positive and can be modeled as a function of relevant predictors via a log-link. Combined with the above, the framework of the proposed model assumes an observation $\mathbf{Y}_t$ at time $t$ is from a Dirichlet distribution with location $\boldsymbol{\theta}$ and scale $\tau$ such that the Dirichlet shape parameter is $\boldsymbol{\alpha} = \tau\boldsymbol{\theta}$. The location parameter is modeled by

$$(2.1) \qquad \boldsymbol{\theta} = \boldsymbol{\eta}/(\boldsymbol{\eta}'\mathbf{1}_p), \quad \text{where } \log(\eta_i) = \beta_{i0} + \beta_{i1}X_1 + \beta_{i2}X_2 + \cdots + \beta_{ik}X_k,$$

and $X_j$, $j = 1, \ldots, k$, are predictor variables with $\beta_{ij}$ as the coefficient on the $j$th predictor for component $i$. Model the scale parameter with

$$(2.2) \qquad \log(\tau) = \gamma_0 + \gamma_1 X_1 + \gamma_2 X_2 + \cdots + \gamma_k X_k,$$

where the $\gamma_j$ terms are the coefficients on the $j$th predictor. It is possible to use different predictor variables in (2.1) and (2.2), and the framework allows for the location/scale to be constant while the other is influenced by covariate terms. It also provides flexibility for one set of the parameters ($\boldsymbol{\theta}$ or $\tau$) to undergo a regime shift while the other stays constant. The predictor variables can model seasonality, trends, or other covariates (ecological drivers).

2.2. *Hidden Markov model.* Let $\mathbf{Y}_t$, $t = 1, \ldots, n$, be an observed composition assumed to follow the Dirichlet distribution. Define $S_t$ to be the latent *state* of the $t$th observation, where $S_t \in \{1, 2, \ldots, m\}$. Further, assume the process controlling the states $S_t$ satisfy the Markov property

$$(2.3) \qquad P(S_{t+1} = j | S_t = i, S_{t-1}, S_{t-2}, \ldots, S_1) = P(S_{t+1} = j | S_t = i) = p_{ij}$$

and that

$$P(\mathbf{Y}_t \in A | S_1, \ldots, S_t = s_t) = P(\mathbf{Y}_t \in A | S_t = s_t).$$

That is, the probability distribution of responses at time $t$ only depends on the underlying state at time $t$. This setup is an $m$ state HMM with a Dirichlet response. Using the formulation in Section 2.1, the latent Markov process effectively determines the regression coefficients in (2.1) and (2.2).

The term $p_{ij}$ in (2.3) is known as the transition probability from state $i$ to state $j$. An $m \times m$ matrix, $\mathbf{P}$ with entry in $i$th row, $j$th column, $p_{ij}$, is known as the probability transition matrix. When a HMM is used for change-point detection, a practitioner may wish to constrain elements of $\mathbf{P}$ to prevent the process from jumping to certain states from others (see Chib (1998)) but an ergodic model may be used in general; see the Supplementary Material for further discussion (Fisher et al. (2022)).

2.3. *Proposed model.* We implement a HMM with the generalized Dirichlet formation presented in Section 2.1. This allows the HMM to detect changes in the underlying location or scale of the distribution. Following Chib (1998), we constrain our transition matrix such that a Markov chain in state $i$ can only jump to state $i + 1$ or remain in state $i$ at the next transition; that is, $p_{ij} = 0$ for all $j \neq i, i + 1$. Figure 3 provides a conceptual schematic of the constrained HMM when three states are present. This constrained HMM, along with the Dirichlet regression models in (2.1) and (2.2), allows us to address the ecological questions: did a considerable shift in phytoplankton phenology occur, and what is the nature of that shift? The Viterbi state assignments (Cappé, Moulines and Rydén (2005)) from the estimated HMM can be used to determine if a change point occurred—a change in Viterbi state indicates a change in the observed distribution.
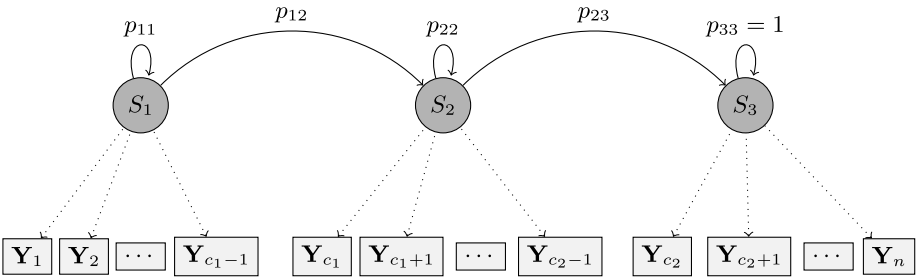
FIG. 3. *Schematic of a constrained three-state hidden Markov model process where change points would be observed in our length n time series at time points $c_1$ and $c_2$.*

2.4. *Model implementation.* This methodology lends itself to both *frequentist* and *Bayesian* implementations (Cappé, Moulines and Rydén (2005)). Maximum likelihood estimation (MLE) for the proposed models is available via the EM-algorithm (Leroux (1992), Lystig and Hughes (2002)), although the constrained HMM may make estimation difficult. The depmixS4 package in R (Visser and Speekenbrink (2010)) provides a set of routines for the practitioner to define their own likelihood and estimate the HMM, even with constraints. Alternatively, HMMs can be implemented in a Bayesian framework (Robert, Celeux and Diebolt (1993)). Due to the constrained nature of our HMM and the small-to-moderate sample size (63 observations on four dimensions), we implement the methods using the latter in the rstan package (Stan Development Team (2018)). Details on the priors of hyperparameters, Bayesian Markov chain Monte Carlo (MCMC) implementation and source code are available in the Supplementary Material (Fisher et al. (2022)).

**3. Simulation study.** The effectiveness of our proposed model is studied through simulation. Data was generated to reasonably mimic the observed seasonal phytoplankton compositions. In the first set of simulations we use dummy variables to model seasonality in the Dirichlet response. Specifically, the design matrix $\mathbf{X}$ is an $n \times 3$ matrix with a first column of all ones corresponding to the *spring*, and the second and third columns contain indicators separating the *summer* and *fall* measurements.

We also implement a variation of the Prabuchandran et al. (2021) Dirichlet likelihood permutation test. For a seasonal response we calculate a different shape parameter for each of the three seasons and perform the permutation within seasonal responses (i.e., spring observations are shuffled with other spring observations), the distribution of the maximum in differences of log-likelihoods was calculated based on 1000 permutations. In the results below, this method is denoted as the *Likelihood Permutation*.

To compare with some other existing methods, we transform the generated compositions via a log-ratio. The transformed data are deseasoned (subtracting seasonal averages), if necessary, and we apply the nonparametric multivariate change implemented in the npcp package in R (see Kojadinovic (2020)), denoted as *Nonparametric Test* below. The nonparametric test implemented in the ecp package (James and Matteson (2014)) is also applied to the transformed data, denoted *ecp Test*. Our simulations were conducted on the *Owens* cluster at the Ohio Supercomputer Center (Ohio Supercomputer Center (2016)). Additional results and details are available in the Supplementary Material (Fisher et al. (2022)).

3.1. *No change data.* Using MLEs for each season of the phytoplankton data in Figure 2, we generate 200 realizations of length $n = 63$ observations, each as four-dimensional compositions, using the parameters in Table 1.

We fit two versions of our proposed model, a single state HMM labeled *Dirichlet Regression* and a two-state model, called *Location & Scale HMM*. To summarize the findings of

TABLE 1
*Parameter values used for simulated datasets when no change point is present*

| Season | $\theta_1$ | $\theta_2$ | $\theta_3$ | $\theta_4$ | $\tau$ |
|---|---|---|---|---|---|
| MLEs for Phytoplankton Data | | | | | |
| Spring | 0.14 | 0.21 | 0.46 | 0.18 | 7 |
| Summer | 0.65 | 0.11 | 0.15 | 0.08 | 11 |
| Fall | 0.55 | 0.14 | 0.22 | 0.09 | 11 |

the models, we considered two simple aggregate metrics to determine whether *convergence* of the posterior simulation routine was achieved and whether the fitted HMM *jumped*. To determine convergence, for each fitted model the average Gelman & Rubin statistics (R-hat from Gelman and Rubin (1992)) value of all model coefficients, the Viterbi states and the log-likelihood value were calculated. If the average R-hat value was less than 1.005, the MCMC was classified as having converged. From there, we randomly sampled 100 of the *converged* fits. We computed the posterior median Viterbi states in each fitted HMM: if the Viterbi state changed, a jump occurred in this fit.

Table 2 reports the results when no change point is present in the data. For the singe state HMM Dirichlet regression ("correct" model), the MCMC algorithm always converged. For the model allowing for a shift in HMM states, the Location & Scale HMM had a 72% rate of convergence; note this model was intentionally misspecified. When the Location & Scale HMM converged, in zero of the 100 randomly selected cases did it jump from state 1 to 2, thus indicating the implemented HMM will not detect nonexistent change points. The nonparametric tests and the Dirichlet likelihood permutation test were calculated on the same 100 randomly selected datasets and report an appropriate number of detected change points (near 5%).

To compare the two fitted Dirichlet-regression based models, we also calculate the approximate leave-one-out (LOO) cross-validation values (Vehtari, Gelman and Gabry (2017)) implemented in the `loo` package (Vehtari et al. (2018)) to assess goodness-of-fit. As expected, the average LOO value across the Dirichlet regression fits is $-383$ with a standard deviation 17, and the LOO for the two-state HMM is $-343$ with standard deviation 19, indicating the single state Dirichlet regression model is the better fit.

3.2. *Single change point.* We now study two scenarios where a change point is present: a seasonal series with a shift in the location parameters (*Location Change Data*) and a seasonal series with a shift in the scale parameters (*Scale Change Data*). The chosen parameters are

TABLE 2
*Proportion of models where the MCMC algorithm converged and proportion of times a change point was detected for each model fit/method when no change point exists*

| Model | MCMC Converged | Change Point Detected |
|---|---|---|
| No Change Data | | |
| Dirichlet Regression | 1.00 | – |
| Location & Scale HMM | 0.72 | 0.00 |
| Nonparametric Test | – | 0.05 |
| ecp Test | – | 0.02 |
| Likelihood Permutation | – | 0.07 |

TABLE 3
*Parameter values used for simulated datasets when a change point occurred*

| Season | State ($s$) | $\theta_1^{(s)}$ | $\theta_2^{(s)}$ | $\theta_3^{(s)}$ | $\theta_4^{(s)}$ | $\tau^{(s)}$ |
|--------|-------------|------------------|------------------|------------------|------------------|--------------|
| Spring | 1 | 0.10 | 0.12 | 0.40 | 0.38 | 13 |
|        | 2 | 0.24 | 0.22 | 0.47 | 0.07 | 1  |
| Summer | 1 | 0.51 | 0.15 | 0.16 | 0.18 | 5  |
|        | 2 | 0.70 | 0.11 | 0.16 | 0.04 | 17 |
| Fall   | 1 | 0.46 | 0.11 | 0.23 | 0.20 | 5  |
|        | 2 | 0.60 | 0.15 | 0.21 | 0.04 | 17 |

based on the MLEs from the phytoplankton data and displayed in Table 3. For the location shift data we use the $\theta_j^{(s)}$ values in Table 3 while keeping $\tau$ constant corresponding to the MLE in Table 1. For the scale shift series we use the respective $\tau^{(s)}$ values in Table 3 while using the MLE values for $\theta_j$ in Table 1. In both scenarios the shift in regime occurs at time point 31.

Table 4 summarizes the results for the two scenarios, comparing the different approaches. There are good rates of convergence for the proposed model in both scenarios, even in the case of single-state Dirichlet regression (misspecified model). The proportion of model fits where a jump occurred is exceptional: 100% of model fits indicated a jump occurred. The two nonparametric tests detect the shift in location while they struggle to detect the shift in scale. The Dirichlet likelihood test modified to handle seasonality provides a perfect rate of detection for both the shift in location and scale scenarios. The LOO results indicate that the Location & Scale HMM is a better fit than the single-state Dirichlet-regression model.

3.3. *Multiple change points.* We also study the methods when 200 realizations of seasonal data with multiple change points are generated. The first 21 observations have a location and scale corresponding to $\theta_j^{(1)}$ and $\tau^{(1)}$ from Table 3, the next 21 observations have the same scale but see a change in location to $\theta_j^{(2)}$ and the remaining 21 observations have a shift in scale $\tau^{(2)}$ while having the same location of $\theta_j^{(2)}$. There are two change points (at time points 21 and 42) with the first experiencing a shift in the location while the second is a shift in scale.

We fit a three-state version of the proposed model to the 200 simulated datasets. For comparison, we implement the multiple change-point algorithm in Prabuchandran et al. (2021)

TABLE 4
*Proportion of models where the MCMC algorithm converged for our approach and proportion of times a change point was detected when one change point exists in the data*

| Model | Location Change Data | | Scale Change Data | |
|-------|-------------------|-------------------|-------------------|-------------------|
|       | MCMC Converged | Change Point Detected | MCMC Converged | Change Point Detected |
| Dirichlet Regression | 1.00 | – | 1.00 | – |
| Location & Scale HMM | 0.98 | 1.00 | 0.99 | 1.00 |
| Nonparametric Test | – | 1.00 | – | 0.14 |
| ecp Test | – | 1.00 | – | 0.22 |
| Likelihood Permutation | – | 1.00 | – | 1.00 |

TABLE 5
*Proportion of models where the MCMC algorithm converged for our model and proportion of times the two change points were detected for simulated data with two change points*

| Model | MCMC Converged | Single Change Point Detected | Two Change Points Detected |
|---|---|---|---|
| Location & Scale HMM | 0.76 | 0.07 | 0.93 |
| Nonparametric Test | – | 0.92 | 0.08 |
| ecp Test | – | 0.94 | 0.06 |
| Likelihood Permutation | – | 0.66 | 0.34 |

and compute the Likelihood permutation test and Nonparametric Test. The ecp Test is designed to segment the data into multiple regimes and is included as well. Table 5 demonstrates the proportion of times our three-state model converged as well as the proportion of times each method detected the different numbers of change points.

All of the approaches consistently detected the first change point, but only the proposed model regularly detects the second change point. Although the Dirichlet likelihood permutation method works well in the single change-point case, it struggles to detect the shift in scale at time point 42, likely due to the limited sample size (only 42 observations, split across three seasons).

Our method has the added feature of modeling the underlying process simultaneous to detecting any shifts. Figure 4 displays the location of change points, based on posterior medians of the Viterbi states in our three-state model; there we see the model accurately detects the true location of the change points, with slightly more variability in detecting a change in scale.

In Figure 5 we display the posterior means of the estimated regression coefficients (the $\beta_{ij}$ and $\gamma_j$ terms in equations (2.1) and (2.2)), along with the theoretical values (gray diamonds) from the simulation parameters. We note that, even with a fairly small sample size (there are only seven observations for each season in each regime), the proposed model estimated the parameters with reasonable accuracy.

3.4. *Covariate influence.* Compared to the methods in the literature, the proposed framework has the added benefit of detecting changes in a model involving covariates. Using the fitted model on the phytoplankton in Section 4.1 as a baseline, we generate a length $n = 63$ seasonal univariate time series, where each *spring* observation is $N(16, \sigma = 1.5)$, *summer*
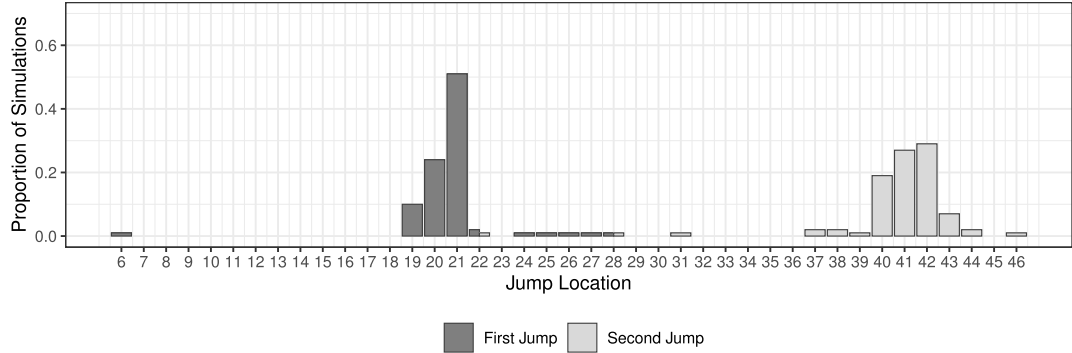


FIG. 4. *Distribution of detected change point locations using our proposed model for simulated data with a shift in location at time point 21 and a change in scale at point 42.*
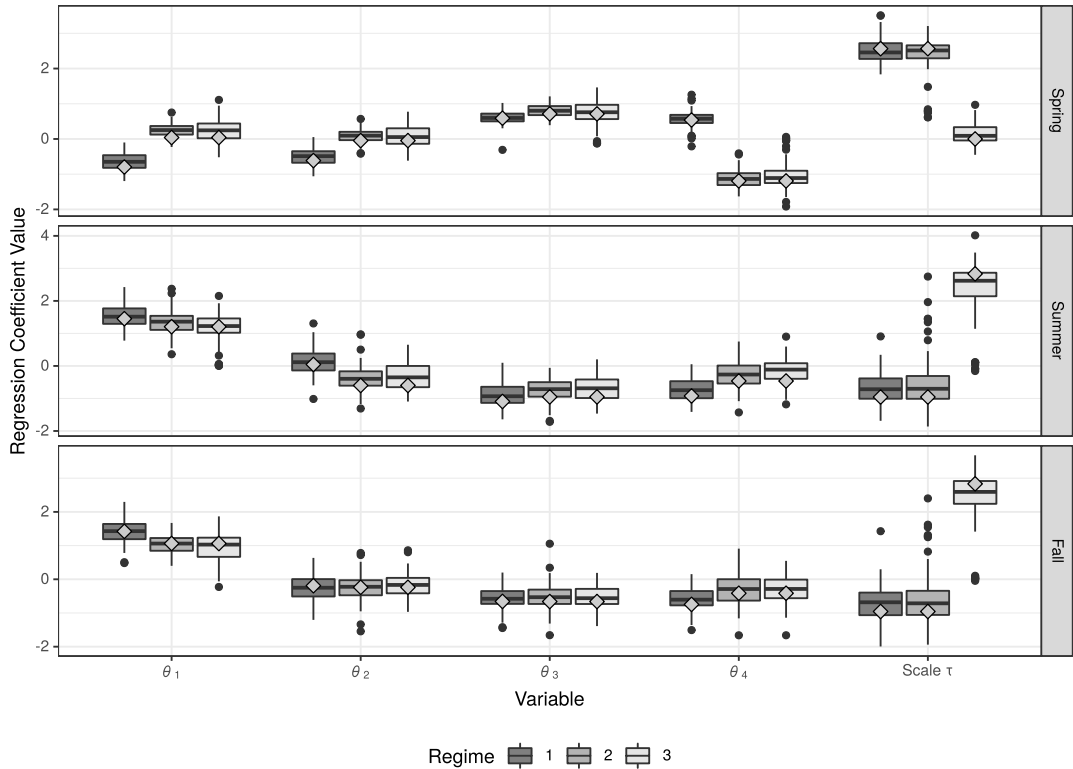
FIG. 5. *Distribution of posterior mean of parameters for the three-state Location & Scale HMM along with the true value (diamond shapes) for 100 randomly selected converged model fits based on data with two change points.*

observations are $N(27, \sigma = 1.5)$ and *fall* observations are $N(21, \sigma = 1.5)$. These covariate terms were used to generate location, and scale parameters using (2.1) and (2.2) with parameters in Table 6, and a compositional response is generated from the corresponding Dirichlet distribution at each time point. The first 30 observations are Regime 1 with the latter 33 in Regime 2.

The proposed Location & Scale HMM with an intercept and the generated covariate as the predictor variable was fit to the simulated data, using the same MCMC parameters as before. In 84% of the 200 simulated datasets, the MCMC algorithm converged. We randomly sampled 100 converged model fits and in 100% of those did the model detect a change point (87 times the change point was detected at time 31, nine times at point 30, and four times at 32). The average (and standard deviation) of the posterior mean of coefficients for each of the selected 100 model fits is in Table 7; there we see estimated coefficients reasonably close to the true values in Table 6.

TABLE 6
*Coefficient values for the covariate model where X is a univariate seasonal time series and the response is generated based on equations (2.1) and (2.2)*

|  | Regime 1 | | | | | Regime 2 | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
|  | $\theta_1$ | $\theta_2$ | $\theta_3$ | $\theta_4$ | $\tau$ | $\theta_1$ | $\theta_2$ | $\theta_3$ | $\theta_4$ | $\tau$ |
| $\beta_{i0}/\gamma_0$ | $-3$ | 0.8 | 1.6 | 0.7 | $-0.2$ | $-1$ | 0.8 | 1.4 | 0.3 | $-0.15$ |
| $\beta_{i1}/\gamma_1$ | 0.2 | $-0.02$ | $-0.20$ | 0.05 | 0.10 | 0.2 | $-0.05$ | $-0.05$ | 0.10 | 0.15 |

| Coefficient | State | $\theta_1$ | $\theta_2$ | $\theta_3$ | $\theta_4$ | $\tau$ |
|---|---|---|---|---|---|---|
| $\beta_{i0}/\gamma_0$ | 1 | −2.71 (0.54) | 0.76 (0.57) | 1.22 (0.63) | 0.74 (0.56) | −0.13 (0.71) |
| | 2 | −1.29 (0.40) | 0.33 (0.60) | 0.89 (0.49) | 0.08 (0.38) | −0.02 (0.67) |
| $\beta_{i1}/\gamma_1$ | 1 | 0.18 (0.06) | −0.03 (0.06) | −0.20 (0.06) | 0.04 (0.06) | 0.10 (0.03) |
| | 2 | 0.16 (0.05) | −0.09 (0.05) | −0.08 (0.06) | 0.05 (0.05) | 0.15 (0.03) |

For comparison, we also fit the version of our model where the design matrix has dummy variables to model seasonality (thus ignoring the covariate influence); that model had a nearly 100% convergence rate and also detected the change point 100% of the time. However, the average LOO goodness-of-fit value for the model with the covariate was −2480 (standard deviation of 396), while the model using seasonal dummy variables was worse fitting with an average LOO of −2429 (standard deviation of 469), thus the model including the covariate influence was a better fit, on average.

3.5. *Additional simulations and conclusions.* The Supplementary Material (Fisher et al. (2022)) provides additional implementation details, the results of additional simulations and variants of the proposed model. In one such simulation we see that by separately modeling the location and scale, we can gain detection power of a change point compared to using the shape.

Simulations indicate the proposed modeling techniques are viable. The HMM can effectively detect shifts in the location and scale parameters of an underlying Dirichlet distribution and works in the multiple change point setting. Since our simulation parameters are based on the observed phytoplankton data, we anticipate our model will work in detecting any regime changes. Based on the provided results (including in Fisher et al. (2022)), the other tests may be helpful in confirming any single change-point shift in location but may struggle to detect multiple changes or a shift in scale.

**4. Analysis of phytoplankton.** We now apply the proposed model to the observed phytoplankton compositions in Figure 2. We fit three versions of our proposed model to the data: a single-state Dirichlet regression model, a two-state HMM, and a three-state HMM. The MCMC algorithm converged (diagnostic plots also available in Fisher et al. (2022)) for all three models with the two- and three-state HMMs showing a change in regimes.

The LOO goodness-of-fit values of the three models are −371.00, −379.25, and −336.72, respectively, indicating that the two-state HMM is the best fit of the three. In fact, the three-state HMM suggest a jump from State 1 to 2, occurring at time points similar to that seen in Figure 6 but in none of the posterior samples did the HMM jump to state 3. The distribution of posterior model coefficients is also similar to those seen in Figures 7 and 8 below, and the posterior $\beta$ and $\gamma$ coefficients in state 3 closely follow the $N(0, 2)$ prior distribution. This suggest that only a single regime shift has occurred in the distribution of phytoplankton.

The two-state HMM suggests a change in all 1000 posterior samples. The posterior distribution of the location of change is in the left panel of Figure 6. The modal peak of that shift is suggested to occur at time point 26 (summer 2000). The right panel of Figure 6 displays the mean posterior Viterbi path demonstrating that, on average, the change point occurs in 2000 and that the years 1999 through 2003 were one of transition for the phytoplankton.

For comparison, we also perform the multiple change-point tests from the literature, each suggests a single change point: the modified nonparametric test (to look for multiple change
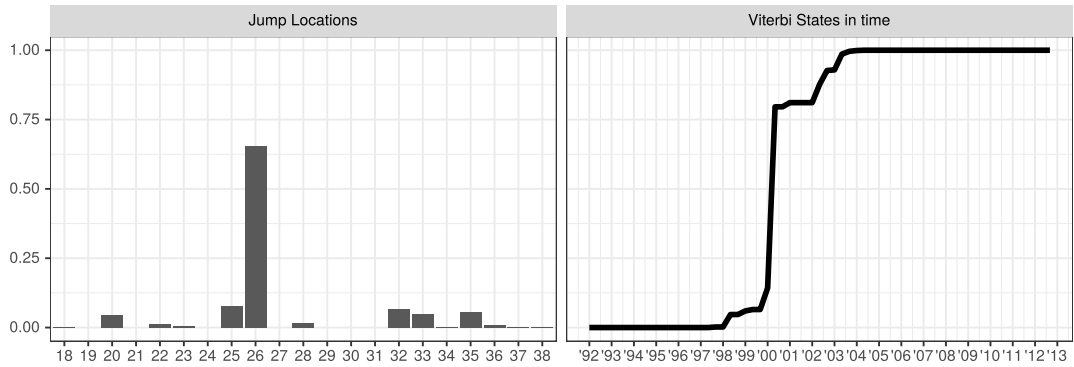
FIG. 6.  *Posterior distribution of jump locations (*left panel*) and mean posterior Viterbi states (*right panel*) for the best fitting Location & Scale HMM for the phytoplankton data.*

points) identifies a change at time point 24 (Fall 1999), the ecp test suggest a change at time 37 (Spring 2004), and the Dirichlet likelihood permutation test at time 34 (Spring 2003). Each of these approaches appear to confirm the detected change, but the proposed method provides more information, as we can see variability associated with any shift in Figure 6.

To gain insight into the nature of the shift suggested by the model, we explore the posterior samples of the estimated Dirichlet parameters. The posterior location estimations ($\theta$), based on the posterior $\beta$ estimates (2.1), were calculated and are displayed in Figure 7. There is a clear decrease in the proportion of diatoms from regime 1 to 2 in all three seasons, with the proportion decreasing to just above zero in the summer and fall seasons. The proportion of



FIG. 7.  *Posterior distribution of estimated phytoplankton location parameters by season and component.*
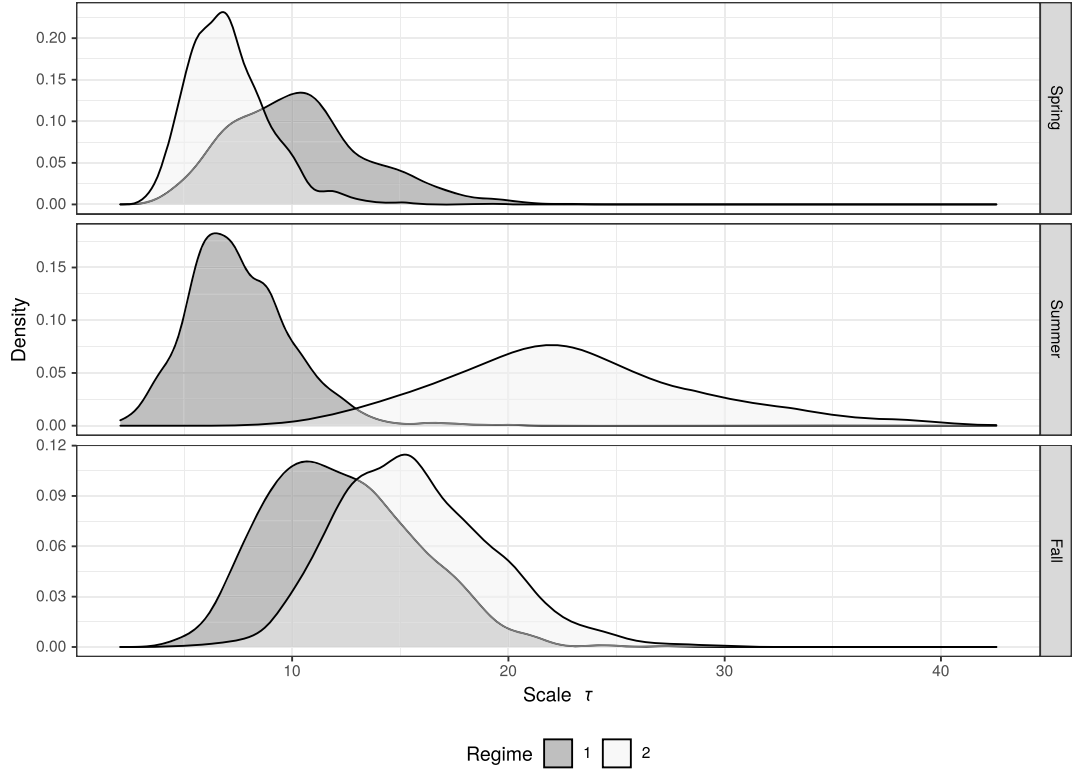
FIG. 8. *Posterior distribution of estimation phytoplankton scale parameters by season.*

blue-green algae increased in the second regime in all three seasons, and there are moderate changes in the proportion of green algae and flagellates (both show a moderate increase in the spring).

The posterior distribution for the scale $\tau$ parameters, calculated from the posterior $\gamma$ estimates (2.2), are displayed in Figure 8. There is fairly strong evidence that the scale in the summer months has increased since the start of the study. Recall that the variance of a component is inversely proportional to the scale term. So, not only do we see an increase in blue-green algae in the summer and a substantial decrease in diatoms, we more consistently see this phenological behavior.

4.1. *Time-varying covariate.* Unlike the other methods studied herein, our framework allows for time-varying covariates in the model, thus allowing for potential *drivers* to explain any observed changes. With the phytoplankton time series of compositions the epilimnion temperature (degrees Celsius), the temperature of the mixed layer at the surface of the lake, where nearly all phytoplankton biomass is contained, is available. The epilimnion temperature is recorded at the time of sampling and aggregated into three measures per year (seasonal averages), as with the phytoplankton. Warming temperatures are known to exacerbate blooms of phytoplankton, especially blue-green algae (Paerl, Otten and Kudela (2018)), so epilimnion temperature may act as an explanatory variable for the phytoplankton phenology.

As expected, the temperature data is highly seasonal (plots available in Fisher et al. (2022)). Thus, using this as a covariate may model the observed seasonality in the phytoplankton compositions. After removing the seasonal effects, the epilimnion temperature time series exhibits some autocorrelation and weak evidence of an underlying trend (trend coefficient: 0.0732 with SE 0.0445 and an AR(1): $\hat{\phi} = 0.2543$, $\text{SE}_{\hat{\phi}} = 0.1208$). The observed in-

crease in epilimnion temperature is similar to many lakes worldwide (O'Reilly et al. (2015)). We note that the trend is not statistically significant (based on historical practice) but is suggestive. Further, the deseasoned data contains 31 negative values and 32 positive values (as expected), but 20 of these negative values occur before time point 36 (99.6% of the posterior paths had a regime shift at or before this point). There is at least suggestive evidence that the epilimnion temperature increased during the period of study and may help explain the observed shift in phytoplankton phenology.

We fit a single state Dirichlet regression model and the two-state HMM with the epilimnion temperature as a covariate. The LOO CV value for the regression model was $-315.7976$ and for the Location & Scale HMM $-283.3801$, suggesting that the regression model is the better fit of the two. From these two models there does not appear to be a regime shift in the underlying Dirichlet regression coefficients when utilizing epilimnion temperature as the only predictor, possibly due to the weak increasing trend in the covariate term. It is worth noting that both LOO CV values for these covariate fitted models are worse than the values reported above. Thus, even though epilimnion temperature is seasonal and appears to have a weak trend, it may not be adequately modeling the the phytoplankton data. Further, the credible intervals for epilimnion temperature regression coefficients for the location cover zero (in all cases) and the sample posterior distributions closely match the prior distributions (see Fisher et al. (2022)).

From the fitted Dirichlet regression model, we calculated the marginal residuals for each phytoplankton group (difference between the observed proportions and the expected proportion) and explored them in time (plot available in Fisher et al. (2022)). From the residuals it is clear that epilimnion temperature is not adequately modeling the seasonality for the blue-green algae proportions or for the flagellates. The expected proportion of green algae appears to be over predicted (mostly negative residuals) and epilimnion temperature does not appear to adequately model the shifts in diatoms. Thus, we can conclude that epilimnion temperature may have some explanatory properties for phytoplankton phenology but is lacking vital information.

The analysis provides evidence the proportion of phytoplankton have changed in time. Furthermore, changes in lake temperature alone do not explain the shift in phenology. Overall, there appears to be a change in the distribution of phytoplankton phenology circa 2000, and the behavior of the taxa appears to follow the distributions presented in Figure 9.
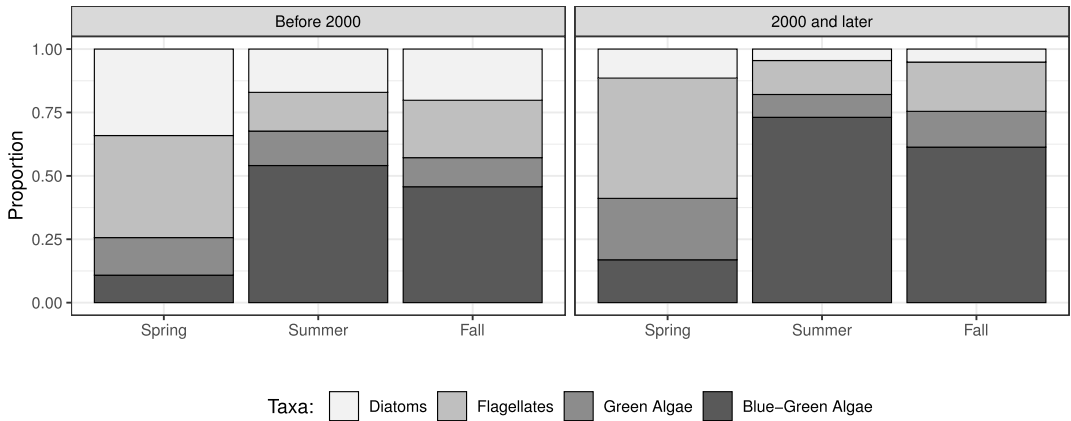


FIG. 9. *Expected compositions of phytoplankton based on the two-state Location & Scale Shift hidden Markov model.*

**5. Conclusions.** In this article we proposed a method for determining if/where change points occur in a compositional time series. The method simultaneously models the effects of the regime shift for the location (expected proportion) and scale (inverse of variance) when assuming a Dirichlet likelihood, the most natural distribution for a composition. Covariate information can be included in the modeling framework. Simulations demonstrate the modeling framework is effective in determining if a shift has occurred, the location of that change, and in estimating the associated parameters of the underlying generalized linear models. In an analysis of 21-years of phytoplankton data (1992 through 2012), we found that there appears to be a shift in the total abundance of phytoplankton circa 2000, based on chlorophyll concentration. Likewise, the proportion of blue-green algae (cyanobacteria) appears to have increased in that time period, while the proportion of diatoms has greatly dissipated. The methods presented here may also have applications in economics (e.g., detect changes in a stock portfolio), geography (detecting changes in land use), and other disciplines.

## SUPPLEMENTARY MATERIAL

`hmmDirichletModel.zip` (DOI: 10.1214/21-AOAS1509SUPPA; .zip). Source code implementing model, simulations and data for analysis.

**Additional results, simulations and data analysis details** (DOI: 10.1214/21-AOAS 1509SUPPB; .pdf). Document outlining further contextual motivation, additional discussion of the proposed model and its implementation, details and simulation results, and further data analysis results.

## REFERENCES

AITCHISON, J. (1982). The statistical analysis of compositional data. *J. Roy. Statist. Soc. Ser. B* **44** 139–177. With discussion. MR0676206

AITCHISON, J. (1985). A general class of distributions on the simplex. *J. Roy. Statist. Soc. Ser. B* **47** 136–146. MR0805071

AITCHISON, J. (1986). *The Statistical Analysis of Compositional Data. Monographs on Statistics and Applied Probability*. CRC Press, London. MR0865647 https://doi.org/10.1007/978-94-009-4109-0

AUE, A. and HORVÁTH, L. (2013). Structural breaks in time series. *J. Time Series Anal.* **34** 1–16. MR3008012 https://doi.org/10.1111/j.1467-9892.2012.00819.x

BARCELÓ-VIDAL, C., AGUILAR, L. and MARTÍN-FERNÁNDEZ, J. A. (2011). Compositional VARIMA Time Series. In *Compositional Data Analysis* 87–103. Wiley, New York.

BARDWELL, L. and FEARNHEAD, P. (2017). Bayesian detection of abnormal segments in multiple time series. *Bayesian Anal.* **12** 193–218. MR3597572 https://doi.org/10.1214/16-BA998

BARRY, D. and HARTIGAN, J. A. (1993). A Bayesian analysis for change point problems. *J. Amer. Statist. Assoc.* **88** 309–319. MR1212493

BINDER, K. E., POURAHMADI, M. and MJELDE, J. W. (2018). The role of temporal dependence in factor selection and forecasting oil prices. *Empir. Econ.* 1–39.

BRUNSDON, T. M. and SMITH, T. M. F. (1998). The time series analysis of compositional data. *J. Off. Stat.* **14** 237–253.

CAPPÉ, O., MOULINES, E. and RYDÉN, T. (2005). *Inference in Hidden Markov Models. Springer Series in Statistics*. Springer, New York. MR2159833

CARLIN, B. P., GELFAND, A. E. and SMITH, A. F. M. (1992). Hierarchical Bayesian analysis of changepoint problems. *J. R. Stat. Soc. Ser. C. Appl. Stat.* **41** 389–405.

CHIB, S. (1998). Estimation and comparison of multiple change-point models. *J. Econometrics* **86** 221–241. MR1649222 https://doi.org/10.1016/S0304-4076(97)00115-2

ERDMAN, C. and EMERSON, J. W. (2008). A fast Bayesian change point analysis for the segmentation of microarray data. *Bioinformatics* **24** 2143–2148.

FEARNHEAD, P. (2006). Exact and efficient Bayesian inference for multiple changepoint problems. *Stat. Comput.* **16** 203–213. MR2227396 https://doi.org/10.1007/s11222-006-8450-8

FEARNHEAD, P. and LIU, Z. (2007). On-line inference for multiple changepoint problems. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **69** 589–605. MR2370070 https://doi.org/10.1111/j.1467-9868.2007.00601.x

FISHER, T. J., ZHANG, J., COLEGATE, S. P. and VANNI, M. J. (2022). Supplement to "Detecting and modeling changes in a time series of proportions." https://doi.org/10.1214/21-AOAS1509SUPPA, https://doi.org/10.1214/21-AOAS1509SUPPB

GELMAN, A. and RUBIN, D. B. (1992). Inference from iterative simulation using multiple sequences. *Statist. Sci.* **7** 457–472.

GRUNWALD, G. K., RAFTERY, A. E. and GUTTORP, P. (1993). Time series of continuous proportions. *J. Roy. Statist. Soc. Ser. B* **55** 103–116.

HAMPTON, S. E., HOLMES, E. E., SCHEEF, L. P., SCHEUERELL, M. D., KATZ, S. L., PENDLETON, D. E. and WARD, E. J. (2013). Quantifying effects of abiotic and biotic drivers on community dynamics with multivariate autoregressive (MAR) models. *Ecology* **94** 2663–2669.

HAYES, N. M. and VANNI, M. J. (2018). Microcystin concentrations can be predicted with phytoplankton biomass and watershed morphology. *Inland Waters* **8** 273–283.

HIJAZI, R. H. and JERNIGAN, R. W. (2009). Modeling compositional data using Dirichlet regression models. *J. Appl. Probab. Stat.* **4** 77–91. MR2668780

HOLMES, M., KOJADINOVIC, I. and QUESSY, J.-F. (2013). Nonparametric tests for change-point detection à la Gombay and Horváth. *J. Multivariate Anal.* **115** 16–32. MR3004542 https://doi.org/10.1016/j.jmva.2012.10.004

JAMES, N. A. and MATTESON, D. S. (2014). Ecp: An R package for nonparametric multiple change point analysis of multivariate data. *Journal of Statistical Software* **62** 1–25.

KANG, S., LIU, G., QI, H. and WANG, M. (2018). Bayesian variance changepoint detection in linear models with symmetric heavy-tailed errors. *Comput. Econ.* **52** 459–477.

KELLY, P. T., GONZÁLEZ, M. J., RENWICK, W. H. and VANNI, M. J. (2018). Increased light availability and nutrient cycling by fish provide resilience against reversing eutrophication in an agriculturally impacted reservoir. *Limnol. Oceanogr.* **63** 2647–2660.

KOJADINOVIC, I. (2020). npcp: Some Nonparametric CUSUM Tests for Change-Point Detection in Possibly Multivariate Observations. R package version 0.2-0.

KOOP, G. and KOROBILIS, D. (2010). Bayesian multivariate time series methods for empirical macroeconomics. *Found Trends Econom.* **3** 267–358.

LEROUX, B. G. (1992). Maximum-likelihood estimation for hidden Markov models. *Stochastic Process. Appl.* **40** 127–143. MR1145463 https://doi.org/10.1016/0304-4149(92)90141-C

LIANG, Z., QIAN, S. S., WU, S., CHEN, H., LIU, Y., YU, Y. and YI, X. (2019). Using Bayesian change point model to enhance understanding of the shifting nutrients-phytoplankton relationship. *Ecol. Model.* **393** 120–126.

LUND, R. and REEVES, J. (2002). Detection of undocumented changepoints: A revision of the two-phase regression model. *J. Climate* **15** 2547–2554.

LUONG, T. M., PERDUCA, V. and NUEL, G. (2012). Hidden Markov model applications in change-point analysis. In *Hidden Markov Models–Applications in Signal*, *Image and Pattern Recognition*.

LYSTIG, T. C. and HUGHES, J. P. (2002). Exact computation of the observed information matrix for hidden Markov models. *J. Comput. Graph. Statist.* **11** 678–689. MR1938450 https://doi.org/10.1198/106186002402

MATTESON, D. S. and JAMES, N. A. (2014b). A nonparametric approach for multiple change point analysis of multivariate data. *J. Amer. Statist. Assoc.* **109** 334–345. MR3180567 https://doi.org/10.1080/01621459.2013.849605

MATTESON, D. S. and TSAY, R. S. (2011). Dynamic orthogonal components for multivariate time series. *J. Amer. Statist. Assoc.* **106** 1450–1463. MR2896848 https://doi.org/10.1198/jasa.2011.tm10616

MILLS, T. C. (2010). Forecasting compositional time series. *Qual. Quant.* **44** 673–690.

O'REILLY, C. M., SHARMA, S., GRAY, D. K., HAMPTON, S. E., READ, J. S., ROWLEY, R. J., SCHNEIDER, P., LENTERS, J. D., MCINTYRE, P. B. et al. (2015). Rapid and highly variable warming of lake surface waters around the globe. *Geophys. Res. Lett.* **42** 10,773–10,781.

OHIO SUPERCOMPUTER CENTER (2016). Owens Supercomputer.

PAERL, H. W., OTTEN, T. G. and KUDELA, R. (2018). Mitigating the expansion of harmful algal blooms across the freshwater-to-marine continuum. *Environ. Sci. Technol.* **52** 5519–5529. PMID: 29656639.

PAWLOWSKY-GLAHN, V. and BUCCIANTI, A., eds. (2011). *Compositional Data Analysis*: *Theory and Applications* Wiley, Chichester. MR2920574 https://doi.org/10.1002/9781119976462

PEARSON, K. (1897). Mathematical contributions to the theory of evolution. – On a form of spurious correlation which may arise when indices are used in the measurement of organs. *Proc. R. Soc. Lond.* **60** 489–498.

PRABUCHANDRAN, K. J., SINGH, N., DAYAMA, P. and PANDIT, V. (2021). Change point detection for compositional multivariate data. *Appl*. *Intell*..

RENWICK, W. H., VANNI, M. J., FISHER, T. J. and MORRIS, E. L. (2018). Stream nitrogen, phosphorus, and sediment concentrations show contrasting long-term trends associated with agricultural change. *J. Environ. Qual*. **47** 1513–1521. https://doi.org/10.2134/jeq2018.04.0162

ROBBINS, M. W., GALLAGHER, C. M. and LUND, R. B. (2016). A general regression changepoint test for time series data. *J. Amer. Statist. Assoc*. **111** 670–683. MR3538696 https://doi.org/10.1080/01621459.2015.1029130

ROBBINS, M., GALLAGHER, C., LUND, R. and AUE, A. (2011a). Mean shift testing in correlated data. *J. Time Series Anal*. **32** 498–511. MR2835683 https://doi.org/10.1111/j.1467-9892.2010.00707.x

ROBBINS, M. W., LUND, R. B., GALLAGHER, C. M. and LU, Q. (2011b). Changepoints in the North Atlantic tropical cyclone record. *J. Amer. Statist. Assoc*. **106** 89–99. MR2816704 https://doi.org/10.1198/jasa.2011.ap10023

ROBERT, C. P., CELEUX, G. and DIEBOLT, J. (1993). Bayesian estimation of hidden Markov chains: A stochastic implementation. *Statist. Probab. Lett*. **16** 77–83. MR1208503 https://doi.org/10.1016/0167-7152(93)90127-5

STAN DEVELOPMENT TEAM (2018). RStan: The R interface to Stan. R package version 2.18.2.

STEPHENS, D. A. (1994). Bayesian retrospective multiple-changepoint identification. *J. Roy. Statist. Soc. Ser. C* **43** 159–178.

TSAY, R. S. (2010). *Analysis of Financial Time Series*, 3rd ed. *Wiley Series in Probability and Statistics*. Wiley, Hoboken, NJ. MR2778591 https://doi.org/10.1002/9780470644560

VEHTARI, A., GELMAN, A. and GABRY, J. (2017). Practical Bayesian model evaluation using leave-one-out cross-validation and WAIC. *Stat. Comput*. **27** 1413–1432. MR3647105 https://doi.org/10.1007/s11222-016-9696-4

VEHTARI, A., GABRY, J., YAO, Y. and GELMAN, A. (2018). loo: Efficient leave-one-out cross-validation and WAIC for Bayesian models. R package version 2.0.0.

VISSER, I. and SPEEKENBRINK, M. (2010). DepmixS4: An R package for hidden Markov models. *J. Stat. Softw*. **36** 1–21.

WEST, M. (2020). Bayesian forecasting of multivariate time series: Scalability, structure uncertainty and decisions. *Ann. Inst. Statist. Math*. **72** 1–31. MR4052647 https://doi.org/10.1007/s10463-019-00741-3