

1  
2 DR. PATRICIA A. SORANNO (Orcid ID : 0000-0003-1668-9271)  
3 DR. KENDRA SPENCE CHERUVELIL (Orcid ID : 0000-0003-1880-2880)  
4 MRS. KATELYN B.S. KING (Orcid ID : 0000-0001-5471-842X)  
5 DR. IAN M MCCULLOUGH (Orcid ID : 0000-0002-6832-674X)  
6 MX. JOSEPH STACHELEK (Orcid ID : 0000-0002-5924-2464)  
7 DR. MERIDITH BARTLEY (Orcid ID : 0000-0001-5896-2948)  
8 MR. NOAH R LOTTIG (Orcid ID : 0000-0003-1599-8144)  
9 DR. TYLER WAGNER (Orcid ID : 0000-0003-1726-016X)

10  
11  
12 Article type : Articles

13  
14  
15 **Journal: Ecological Applications**

16 **Manuscript type: Article**

17  
18  
19 **Running head:** Ecological prediction at macroscales

20  
21  
22 **Ecological prediction at macroscales using big data: Does sampling design matter?**

23  
24 **AUTHORS:** Patricia A. Soranno<sup>1</sup> and Kendra Spence Cheruvellil<sup>1,2,10</sup>, Boyang Liu<sup>3</sup>, Qi Wang<sup>3</sup>,  
25 Pang-Ning Tan<sup>3</sup>, Jiayu Zhou<sup>3</sup>, Katelyn B.S. King<sup>1</sup>, Ian M. McCullough<sup>1</sup>, Joseph Stachelek<sup>1</sup>,  
26 Meredith Bartley<sup>4</sup>, Christopher T. Filstrup<sup>5</sup>, Ephraim M. Hanks<sup>4</sup>, Jean-Francois Lapierre<sup>6</sup>, Noah  
27 R. Lottig<sup>7</sup>, Erin M. Schliep<sup>8</sup>, Tyler Wagner<sup>9</sup>, Katherine E. Webster<sup>1</sup>

28

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1002/EAP.2123](https://doi.org/10.1002/EAP.2123)

This article is protected by copyright. All rights reserved

29 <sup>1</sup>Department of Fisheries and Wildlife, Michigan State University, 480 Wilson Road, East  
30 Lansing, Michigan 48824, USA

31 <sup>2</sup>Lyman Briggs College, Michigan State University, 919 East Shaw Lane, East Lansing,  
32 Michigan 48825, USA

33 <sup>3</sup>Department of Computer Science and Engineering, Michigan State University, 428 South Shaw  
34 Lane, East Lansing, Michigan 48824, USA

35 <sup>4</sup>Department of Statistics, The Pennsylvania State University, 324 Thomas Building, University  
36 Park,  
37 Pennsylvania 16802, USA

38 <sup>5</sup>Natural Resources Research Institute, University of Minnesota Duluth, 5013 Miller Trunk  
39 Highway, Duluth, Minnesota 55811, USA

40 <sup>6</sup>Sciences Biologiques, Universite de Montreal, Pavillon Marie-Victorin, CP 6128, succursale  
41 Centre-Ville, Montreal, Quebec, H3C 3J7, CANADA

42 <sup>7</sup>Center for Limnology Trout Lake Station, University of Wisconsin Madison, Boulder Junction,  
43 Wisconsin, 54512, USA

44 <sup>8</sup>Department of Statistics, University of Missouri, 146 Middlebush Hall, Columbia, Missouri  
45 65211, USA

46 <sup>9</sup>U.S. Geological Survey, Pennsylvania Cooperative Fish and Wildlife Research Unit,  
47 Pennsylvania State University, Forest Resources Building, University Park, Pennsylvania  
48 16802, USA

49 <sup>10</sup>Corresponding author. E-mail: [ksc@msu.edu](mailto:ksc@msu.edu)

50  
51  
52 Manuscript received 9 July 2019; revised 13 December 2019; accepted 6 January 2020; final  
53 version received 9 February 2020.

54  
55 **ABSTRACT**  
56 Although ecosystems respond to global change at regional to continental scales (i.e.,  
57 macroscales), model predictions of ecosystem responses often rely on data from targeted  
58 monitoring of a small proportion of sampled ecosystems within a particular geographic area. In  
59 this study, we examined how the sampling strategy used to collect data for such models  
60 influences predictive performance. We subsampled a large and spatially-extensive dataset to

61 investigate how macroscale sampling strategy affects prediction of ecosystem characteristics in  
62 6,784 lakes across a 1.8 million km<sup>2</sup> area. We estimated model predictive performance for  
63 different subsets of the dataset to mimic three common sampling strategies for collecting  
64 observations of ecosystem characteristics: random sampling design, stratified random sampling  
65 design, and targeted sampling. We found that sampling strategy influenced model predictive  
66 performance such that (1) stratified random sampling designs did not improve predictive  
67 performance compared to simple random sampling designs and (2) although one of the scenarios  
68 that mimicked targeted (non-random) sampling had the poorest performing predictive models,  
69 the other targeted sampling scenarios resulted in models with similar predictive performance to  
70 that of the random sampling scenarios. Our results suggest that although potential biases in  
71 datasets from some forms of targeted sampling may limit predictive performance, compiling  
72 existing spatially-extensive datasets can result in models with good predictive performance that  
73 may inform a wide range of science questions and policy goals related to global change.

74

75 **KEYWORDS:** extrapolation, interpolation, lakes, prediction, sampling design, macroscale,  
76 data-intensive ecology, monitoring, sampling, ecological context

77

## 78 **INTRODUCTION**

79 Scientific evidence from focused monitoring efforts has been used since the 1990's to  
80 inform environmental policy in response to broad-scale environmental stressors such as acid rain  
81 and lake eutrophication (Olsen et al. 1999), and there has been much interest in knowing how  
82 different strategies used to select sample ecosystems may affect inference (e.g., Janousek et al.  
83 2019). Previous work has been conducted primarily at local to regional scales, often focusing on  
84 geographic areas containing the most sensitive ecosystems. In recent years, there has been a  
85 growing recognition of the need to predict ecosystem responses to global change over broader  
86 spatial extents that encompass scales from regions to continents (Miller et al. 2004, Dietze et al.  
87 2018, Peters et al. 2018; hereafter referred to as macroscales *sensu* Heffernan et al. 2014). To  
88 date, it is unknown how sampling design affects our ability to understand and predict states and  
89 relationships in unsampled ecosystems at macroscales.

90 Prediction at macroscales is complicated because it requires integration of the multi-  
91 scaled spatial variation that underlies temporal responses to drivers of global change. Because  
92 spatial heterogeneity can be large and can exceed temporal variation (Soranno et al. 2019), it is a

93 critical component to be accounted for when predicting ecosystem states and relationships at the  
94 macroscale. Further, prediction accuracy is strongly influenced by the spatial variation of the  
95 data used to generate models, which means that the strategy used to select sample ecosystems  
96 plays a large role in predictive modeling success (Thompson 2012).

97 There are two main ways to acquire data for predictive models at the macroscale –  
98 coordinated national monitoring programs and compilations of more localized (e.g., local or  
99 regional) and disparate datasets. Examples of the first approach include the U.S. Environmental  
100 Protection Agency’s National Lakes Assessment program that samples approximately 1,000  
101 lakes every five years, comprising ~1 % of lakes  $\geq 1$  ha (U.S. Environmental Protection Agency  
102 Office of Wetlands, Oceans and Watersheds Office of Research and Development 2017).  
103 Similarly, the U.S.D.A. Forest Service’s Forest Health Monitoring Program samples  
104 approximately 12,500-25,000 plots annually, comprising 10-20% of all forest plots (Smith  
105 2002). A recent example of the second approach is a macroscale dataset of lake observations  
106 created by compiling almost 90 disparate local and regional datasets across 17 U.S. states  
107 resulting in approximately 12,000 lakes with at least one observation, comprising 24% of lakes  $\geq$   
108 4 ha (Soranno et al. 2017). In both approaches, a small proportion of ecosystems is sampled and  
109 the knowledge gleaned from them is consequently applied to unsampled ecosystems.

110 Various strategies have been used to select ecosystems for sampling in macroscale  
111 monitoring programs in the past, each with their strengths and weaknesses in terms of resources  
112 required, potential biases introduced, and predictive power (Urquhart et al. 1998, Olsen et al.  
113 1999, Thompson 2012, Sauer et al. 2013). At the macroscale, sample ecosystems are rarely  
114 selected using a simple random design but are sometimes selected using a stratified random  
115 design. There has also been a long history of sampling ecosystems for purposes such as  
116 ecosystem management without using a probabilistic sampling design that allows representation  
117 of the entire population. In these cases, targeted sampling is conducted for subsets of ecosystems  
118 or landscapes that are of interest, such as regions that are of high conservation interest or  
119 ecosystems that are at high-risk of human perturbation (i.e., observational studies where there is  
120 little to no control over which ecosystems are sampled; Thompson 2012). None of these  
121 strategies result in a dataset that is a perfectly representative sample of the entire population,  
122 particularly when using the sample data for prediction of unsampled ecosystems. In practice, the  
123 majority of existing macroscale datasets are likely to be biased in different ways, with some  
124 datasets over- or others undersampling particular ecosystems or those with particular

125 characteristics (e.g., Webb et al. 2013, Stanley et al. 2019, Zhao et al. 2019). For example, when  
126 multiple disparate datasets are compiled, the resulting datasets include data from a mixture of  
127 probabilistic sampling designs and targeted sampling efforts, the effects of which can only be  
128 quantified after the database has been created (e.g., GBIF, LAGOS-NE; Gaiji et al. 2013,  
129 Soranno et al. 2017).

130         When building a predictive model, it is a common practice to train the model using a  
131 subset of the dataset (training data) and then test the model using data that were withheld (out-of-  
132 sample or test data; Lohr 2019). A fundamental assumption behind most predictive empirical  
133 models is that the training and test data are generated from the same distributions (i.e.,  
134 predictions made within the model space). Thus, the resulting predictions are thought of as  
135 interpolations. However, if the training and test data are from different distributions, then there is  
136 no guarantee that the model fitted to the training data will perform well on the test data (i.e.,  
137 predictions made outside the model space). For example, predictions at unsampled ecosystems  
138 with predictor variables that exceed the range of predictors in the training data and/or comprise a  
139 novel combination of predictors may be unreliable and are commonly referred to as  
140 extrapolations (Conn et al. 2015). Encountering such novel settings may occur often in  
141 macroscale studies due to the broad spatial extent associated with them and the large gradients  
142 that exist at these extents for the many characteristics that make up an ecosystem's ecological  
143 context (e.g., land use/cover, geology, climate). Therefore, it is critical to assess how various  
144 sampling strategies with different purposes may introduce biases that affect distributions of  
145 training and test datasets and could change interpolations to extrapolations, thus influencing  
146 model predictive performance.

147         We used a large database compiled from local and regional disparate datasets to ask:  
148 what is the effect of sampling design on predictive models of ecosystem characteristics in  
149 unsampled ecosystems at the macroscale? We used 4,253-6,784 observations of lake nutrients  
150 and productivity from a dataset of 51,101 lakes and their ecological contexts within a spatial  
151 extent of 1,778,100 km<sup>2</sup> in the northeastern and midwestern U.S. to answer this question  
152 (Soranno et al. 2015, 2017). Although this database has its own inherent biases (e.g.,  
153 undersampling of small lakes; Stanley et al. 2019), it includes a wide range of lake types with  
154 large gradients of ecosystem characteristics located across many regions with large gradients of  
155 ecological contexts. Therefore, it is an ideal database to create subsets of data that represent

156 known degrees of bias in order to quantify the effects of sampling design on predictive models of  
157 ecosystem characteristics.

158 We developed scenarios (described below) that mimic three common strategies used for  
159 collecting observations on ecosystem characteristics at macroscales: random sampling design,  
160 stratified random sampling design, and targeted sampling. We used three measures of lake  
161 ecosystem characteristics, total phosphorus, total nitrogen, and chlorophyll *a*, to compare the  
162 predictive performance of models across these scenarios and strategies. We expected models to  
163 have highest predictive performance in cases of assumed interpolation and lowest in cases of  
164 assumed extrapolation (Conn et al. 2015). We also expected stratified random designs to increase  
165 predictive performance of the interpolation scenarios because the strata are chosen based on  
166 underlying ecological processes that are more likely to be related to spatial variation than strictly  
167 random sampling. Therefore, we expected predictive performance to be highest when using the  
168 stratified random designs, moderate when using the random designs, and lowest when using the  
169 targeted sampling. We also expected better predictive performance when using a relatively large  
170 proportion of lakes to train or build the predictive model. Finally, we expected nutrients, which  
171 are directly linked with landscape context variables, to be better-predicted than lake productivity.  
172 Our results will inform the design of macroscale ecosystem assessments, lead to more robust  
173 understanding of macroscale variation among ecosystems, and result in better predictions of  
174 unsampled ecosystems.

175

## 176 **Conceptualizing the effect of sampling design on predictive models of unsampled** 177 **ecosystems**

178 We created seven scenarios that fall within one of the three common sampling strategies  
179 employed in macroscale studies, the data from which are used to develop models used to predict  
180 at unsampled ecosystems. Figure 1 depicts these strategies as columns with multiple scenarios  
181 under each strategy labeled (a-g) and with training data in orange and test data in blue. The  
182 scenarios depicted in the left panel of Figure 1 (a-b) illustrate the rare cases when ecosystems are  
183 selected at random at the macroscale, called *random sampling design*. Figure 1a depicts the best-  
184 case scenario whereby a large proportion of the data are used to train the model and a small  
185 proportion of data are used to test the model. We use this scenario as a predictive baseline to  
186 compare with the other scenarios that have smaller training datasets since having large datasets  
187 to build predictive models is extremely rare in ecology and those that are available often have

188 been compiled from multiple (non-random) sources that are question- or problem-driven. Figure  
189 1b shows the more common scenario in which a smaller dataset is available for model training  
190 and the test dataset is larger. If the sample size of the training dataset is sufficiently large, model  
191 predictions in these cases are assumed to be within the model space, resulting in *interpolation*.

192 A second set of scenarios demonstrate stratified random sampling designs that are  
193 commonly used in macroscale ecosystem assessments (Figure 1c–d). In these cases, factors that  
194 are thought to be important for driving ecosystem processes and patterns are used to first stratify  
195 the entire population of ecosystems and then ecosystems are randomly selected within each  
196 stratum. Figure 1c–d represents two common cases of *stratified random sampling design*. Figure  
197 1c depicts sample selection using strata based on ecosystem type and Figure 1d depicts selection  
198 using strata based on spatial location of the ecosystem, i.e. region. Species presence and/or  
199 abundance are commonly estimated with stratified sampling designs whereby the landscape is  
200 stratified by ecologically important characteristics (e.g., moose surveys across vegetation types  
201 or high/low quality habitat or fish surveys in lakes stratified by depth and area; Ver Hoef 2008,  
202 Rask et al. 2010). Stratified random designs assume that the feature(s) used to define strata are  
203 ecologically relevant for the response variables being considered by the study (i.e., ecosystem  
204 type or regions drive variation among ecosystems). Therefore, as long as the sample size of the  
205 training dataset is sufficiently large, predictions for unsampled ecosystems are assumed to be  
206 *interpolations*.

207 The third common sampling strategy is targeted sampling that happens when assessments  
208 are question- or problem-driven. In these cases, particular ecosystem types or regions are  
209 targeted for sampling in order to answer specific questions or to assess specific populations of  
210 ecosystems. Two examples of this design are giant sequoia trees sampled to reconstruct regional  
211 fire histories (Swetnam 1993) and lakes in the U.S. sampled as part of the National  
212 Eutrophication Survey to study causes of eutrophication (US Environmental Protection Agency  
213 1975). Figure 1e–g depict examples of *targeted sampling* that result in the training datasets being  
214 based on particular ecosystem types (Figure 1e), regions (Figure 1f), or regions with particular  
215 land uses (Figure 1g). In such cases of targeted sampling, when the data are used in predictive  
216 models of unsampled ecosystem types or regions, we assume that these ecosystems are not  
217 representative of all ecosystems. Therefore, we assume that these may represent cases of  
218 *extrapolation*.

219

## 220 **METHODS**

### 221 **Study site and dataset**

222 We used the LAGOS-NE database that spans the lake-rich regions of the northeastern  
223 and midwestern U.S. (Soranno et al. 2015, 2017) and includes 4,253 – 6,784 lakes depending on  
224 the response variable (from a total population of 51,101 lakes  $\geq 4$  ha). The study lakes include  
225 both shallow and deep lakes (interquartile range of maximum depth = 4.6-13.7 m), natural lakes  
226 and reservoirs, and lakes with watersheds that are entirely forested to entirely surrounded by  
227 agricultural land use. The lakes in this database cover broad gradients in climate, geology, land  
228 use/cover, hydrology, and topography. LAGOS-NE-GEO v1.05 includes lake, local, and regional  
229 ecological context (Soranno and Cheruvilil 2017a) and LAGOS-NE-LIMNO v1.087.1 includes at  
230 least one *in situ* observations of lake water quality for 10,173 lakes (Soranno and Cheruvilil  
231 2017b). These lakes are nested within 65 regions defined by the level 4 hydrologic units  
232 regionalization (Seaber et al. 1987; hereafter referred to as regions and HU4s; Figure 1). These  
233 regions with an average area of 43,500 km<sup>2</sup>, have been shown to account for regional variation in  
234 nutrients and productivity of this lake population (Cheruvilil et al. 2013). Data and code are  
235 available (see Data Availability).

### 236 **Lake response variables**

237 We analyzed three ecosystem characteristics of lakes that represent major nutrients and  
238 primary productivity – total phosphorus (TP), total nitrogen (TN), and chlorophyll *a* (CHL).  
239 These variables are routinely measured by a wide range of academic, governmental, and non-  
240 governmental programs to assess water quality (Poisson et al. 2019). We selected lakes and  
241 observations using the following criteria. Lake nutrient and productivity observations were  
242 selected during the time of peak production in these lakes (i.e., the summer stratified period of 15  
243 June through 15 September) during the years 1980 to 2011. For lakes with multiple observations  
244 within a summer or across multiple years, we selected a single sample that contained the most  
245 response variables. The resulting data came from lakes ranging from very nutrient-poor and low  
246 productivity systems to very nutrient-rich and high productivity systems that are distributed  
247 across our study area (Table 1; Figure 1).

### 248 **Local and regional ecological context predictor variables**

249 We selected 18 predictor variables *a priori* that are consistently related to lake nutrients  
250 and productivity (Table 2; Read et al. 2015, Collins et al. 2017, Lapierre et al. 2018, Soranno et  
251 al. 2019). At the local scale, we included six lake-specific characteristics – lake connectivity type



252 (defined as either lakes that have either no stream connections or only outflowing stream  
253 connections (Isolated), lakes with inflowing and outflowing stream connections (DR\_Stream), or  
254 lakes with connections to upstream lakes (DR\_LakeStream); lake water clarity (as measured by  
255 Secchi disk depth); maximum lake depth; lake complexity (a metric of lake shape that measures  
256 the deviation of the shoreline from a circular shape); and, lake elevation. Lake water clarity was  
257 included because it is available for nearly all lakes in our study sample (Figure 2) and model  
258 predictions are more accurate when they are conditional on water clarity (Wagner et al. 2020,  
259 Wagner and Schliep 2018). We also included five watershed-specific characteristics for the area  
260 of land draining directly into the lake as well as the area that drains into upstream-connected  
261 streams and lakes <10 ha (i.e., the inter-lake watershed; Soranno et al. 2017) – watershed  
262 wetland cover; watershed complexity (a metric of watershed shape that measures the deviation of  
263 the watershed boundary from a circular shape); watershed to lake area ratio; watershed stream  
264 density; watershed forest cover; and, watershed road density. Finally, seven regional-scale  
265 characteristics calculated for each HU4 were included in models – mean percent baseflow (an  
266 index of regional groundwater contribution); mean runoff; percent agricultural land use; mean  
267 annual temperature; mean annual precipitation; mean total nitrogen deposition in 1990; and, the  
268 difference in mean total nitrogen deposition from 1990 to 2010. Details on the data sources for  
269 these variables are provided in Soranno et al. (2017).

## 270 **Macroscale sampling scenarios**

271 We created seven sampling scenarios that mimic common approaches used for collecting  
272 observations on ecosystem characteristics at the macroscale (Figure 1a-g). In these scenarios, we  
273 assumed that the population of LAGOS-NE lakes with TP ( $n = 5,896$ ), TN ( $n = 4,253$ ), or CHL  
274 ( $n = 6,784$ ) represent the census population (but see Stanley et al. 2019), and that the training and  
275 test data were subsets of this population. We fitted models to each of these seven scenario  
276 datasets and compared predictive performance (see below for details) for modeling the state of  
277 ‘unsampled’ ecosystems in the test dataset.

278 *Random sampling designs:* In these two scenarios, we used a random sampling design  
279 and examined the effect of sample size on predictive model performance (Figure 1a-b). First, we  
280 created a scenario that represents an analytical predictive baseline with a large training dataset of  
281 75% of sampled lakes (Figure 1a). As a contrast, we created a scenario that uses a small training  
282 dataset of 25% of sampled lakes (Figure 1b).

283 *Stratified random sampling designs:* In these two scenarios, we stratified the sampling  
284 based on ecological context measured at either the local scale (based on lake type) or the regional  
285 scale (based on the region membership of each lake) (Figure 1c-d). For the lake type strata, we  
286 created four clusters of lakes based on watershed and regional landscape context characteristics  
287 (Table 2) and using hierarchical clustering using Ward's method (Ward 1963). Cluster 1 was  
288 characterized by lakes in regions with above average number of and extent of upstream lakes.  
289 For the remaining three clusters that had below average regional upstream lake connectivity,  
290 lakes were characterized by either high stream density in the watershed (cluster 3), high percent  
291 of wetlands in the watershed and around the lake perimeter (cluster 4,) or by both low stream  
292 density and low wetland percent in the watershed (cluster 2). For both stratified random  
293 scenarios, we selected 25% of lakes within each strata (lake type or region) to build the  
294 predictive models, and then predicted the values for the remaining 75% of lake ecosystems as we  
295 did for the random sampling design scenario that had a small training dataset (described above).

296 *Targeted sampling:* We created three targeted sampling scenarios by selecting lakes of  
297 particular types, particular regions, or particular types of regions. First, using the above four lake  
298 type clusters, we selected all lakes in two of the four clusters to form the training dataset, and  
299 tested the model on the lakes in the remaining two of the four clusters. Second, we selected all  
300 lakes in half of the regions to form the training dataset and tested the model on lakes in the  
301 remaining half of the regions. For these two targeted scenarios, we split the sampled lake data  
302 approximately 50:50 and randomly selected the lake type clusters or regions for training the  
303 models. For the third targeted scenario, we deliberately selected half of the regions with the  
304 lowest proportion of agricultural land to form the training dataset and tested the model on lakes  
305 in the remaining regions with the highest proportion of agricultural land. However, because lakes  
306 are not distributed equally across regions, the number of lakes was not 50:50 in the training:test  
307 datasets for this scenario. The high-agriculture regions contain only 25% of the sampled lakes in  
308 the study area, whereas the low-agricultural regions are very lake-rich and contain 75% of the  
309 sampled lakes in the study area.

### 310 **Predictive models of ecosystem characteristics**

311 We used random forest models (Breiman 2001, Liaw and Wiener 2002) to predict each of  
312 the three response variables (TP, TN, CHL) based on the 18 local (lake-specific and watershed)  
313 and regional predictor variables described above that are related to lake nutrients and  
314 productivity in the LAGOS-NE lakes (Table 1, 2). Random forest is an ensemble learning

315 method that generates its prediction by averaging the outputs produced by a set of regression  
316 trees; and, each regression tree is created via bootstrapping by applying sampling with  
317 replacement on the training data (Breiman 2001, Zhou 2012). Although there are no  
318 distributional assumptions for random forests, the algorithm determines the best model based on  
319 squared error between predictions and true, out of sample, data (Breiman 2001).

320 We log-transformed the response data after adding 0.1 to the values to down-weight  
321 errors on lakes with large data values so that our error terms are closer to percent error than to  
322 absolute error. For predictor variables, there were a few cases of missing values (1.97% of  
323 values). Those values were imputed with the mean value for that variable so that all observations  
324 could be used in the random forest models. The predictor variables were standardized by  
325 subtracting the mean and dividing by the standard deviation.

326 After these pre-processing steps, the dataset was split into training and test datasets based  
327 on the seven scenarios depicted and described above (Figure 1). To minimize the likelihood of  
328 chance selection affecting modeling results, we randomly split the dataset into training and  
329 testing datasets 10 times for each scenario possible (four of the seven scenarios). For the two  
330 scenarios that used four lake types, we could only create six sets of training and testing datasets  
331 (all possible combinations of four types). For the targeted sampling scenario that used regional  
332 land use, only one train/test dataset could be created.

333 We then used the random forest method, building 189 total independent models, one for  
334 each combination of response variable (3), sampling design scenario (7), and train/test dataset  
335 combination (1, 6, or 10 as described above). Random forest has several hyperparameters that  
336 need to be tuned with the training data, including maximum tree depth and number of trees. We  
337 conducted a grid search, with both of these hyperparameters allowed to range from 50 to 200.  
338 We performed 5-fold cross validation (Stone 1974) on the training data to determine the optimal  
339 hyperparameter setting. Specifically, we iteratively reserved four of the five folds for model  
340 building and used the remaining fifth fold as a validation set to select the best hyperparameters.  
341 We then re-trained the random forest model on the entire training set using the best  
342 hyperparameter values and applied the resulting model to the test dataset to predict the response  
343 variables. We trained our random forest model using the Python scikit-learn RandomForest  
344 package with Gini impurity as the splitting criterion of the tree (Pedregosa et al. 2011).

345 *Predictive performance:* We quantified model predictive performance three ways for  
346 each of the 189 independent models to compare the effect of sampling scenarios on model

347 performance. First, we calculated the root mean squared error (RMSE), which is a measure of  
348 average prediction error that is in the units of the log-transformed response variable. Second, we  
349 calculated the median relative absolute error (MRAE;  $\varepsilon = \text{median}(\frac{|\hat{y} - y|}{y})$ ), which is a unitless  
350 measure of relative error that can be useful for comparing model performance across response  
351 variables. Third, we calculated the predictive  $R^2$ , which is a bounded measure of model relative  
352 accuracy whereby 0 indicates that model prediction is no better than using the mean value of the  
353 response variable and 1 indicates perfectly accurate model prediction. For the five scenarios with  
354 multiple train/test dataset combinations, we calculated average predictive performance and  
355 corresponding standard error over the multiple train/test dataset combinations.

356

## 357 RESULTS

358 The predictive models accounted for 34-63% of the variation in lake nutrients and  
359 productivity across the seven scenarios that mimic three common macroscale sampling strategies  
360 (Figure 3A). In general,  $R^2$  values decreased from larger to smaller training datasets, from  
361 random sampling design (stratified or not) to targeted sampling, and from TP to TN and to CHL.  
362  $R^2$  was  $> 0.5$  for all of the response variables and sampling scenarios, except for when modeling  
363 TN using the regional land use targeted sampling scenario (Figure 3A (g)). The predictor  
364 variables that accounted for most of the variation in responses were lake and watershed  
365 landscape characteristics such as lake maximum depth, watershed percent forest, and water  
366 clarity (Appendix S1: Table S1-S3).

367 The two random sampling design scenarios are (a) and (b) in Figure 3. The scenario that  
368 used 75% of lakes to build the model (a) resulted in predictions of nutrients and productivity  
369 with the lowest error as measured by RMSE and MRAE (Figure 3B-C). Although we considered  
370 this scenario to be somewhat unrealistic in practice due to the large sample size (e.g.,  $n = 4,422$   
371 for TP), when we decreased that sample size to 25% of sampled lakes (e.g.,  $n = 1,474$  for TP;  
372 (b)), the effect on predictive performance was negligible (change in RMSE of 0.02-0.03; Figure  
373 3B, Table 3).

374 The two stratified random sampling design scenarios are (c) and (d) in Figure 3. When  
375 comparing the simple random sampling design scenarios with the smaller training dataset (b) to  
376 these two stratified random sampling designs, we found small differences in predictive  
377 performance (Figure 3, Table 3). In fact, the differences among these three random sampling  
378 design scenarios (stratified or not) were nonexistent to negligible. Therefore, the three assumed

379 interpolation scenarios with smaller training datasets (b – d) were similarly able to predict lake  
380 nutrients and productivity.

381 The three scenarios that represent targeted sampling are (e) - (g) in Figure 3. Targeted  
382 sampling based on lake type (e) or region (f), resulted in slightly lower predictive performance  
383 and higher variation across simulated datasets compared to the random sampling design scenario  
384 that uses the smaller training dataset (b) (Figure 3). However, the scenario that mimicked  
385 targeted sampling of regions with high agriculture (g) resulted in the poorest performance of any  
386 scenario, particularly for TN. This poor performance is likely due to this scenario being a case of  
387 extrapolation as demonstrated by differences in the distributions of the response variables and  
388 important predictor variables between the training and testing datasets for (g) that was not  
389 apparent for the other scenarios (Figure 4).

390

## 391 **DISCUSSION**

392 We studied 6,784 lakes across a spatial extent of 1.8 million km<sup>2</sup> to understand how  
393 different sampling strategies may affect model predictions of commonly measured ecosystem  
394 characteristics in unsampled ecosystems at macroscales. We found that although the sampling  
395 strategy used is likely to influence model predictive performance, the differences may not always  
396 be as large or as expected based solely on sample sizes and whether the strategy results in  
397 interpolation or extrapolation. We have two specific take home messages from this research.  
398 First, sampling designs based on two commonly used stratified random approaches (i.e., by  
399 region or by ecosystem type) did not result in better predictions of lake nutrients and productivity  
400 compared to a simple random sampling design, suggesting that at the macroscale, stratified  
401 random sampling designs may not *always* be better than simple random sampling designs.  
402 Second, models trained with data from targeted sampling were not always the poorest  
403 performing models. However, the predictive performance varied across the three targeted  
404 sampling scenarios and three response variables. This fact suggests that data from some targeted  
405 sampling may result in extrapolation and poor model performance, and thus should be examined  
406 for potential biases before use. Below, we discuss the effects of sampling strategies on predictive  
407 model performance and interpret these effects within the context of LAGOS-NE, the database  
408 that was used to create the seven sampling scenarios. Then, we discuss the implications of our  
409 results for optimizing macroscale sampling designs.

410

## 411 **Effects of sampling strategies on predictive performance**

412 We anticipated that random and stratified random sampling designs would outperform  
413 targeted sampling. This expectation was based mainly on the assumption that targeted sampling  
414 designs would result in the training and test data having different distributions, meaning  
415 predictions would be made outside of the model space (i.e., extrapolation). However, our results  
416 demonstrated that this assumption does not always hold true. For example, the distributions of  
417 the training and testing datasets for the response and predictor variables were very similar for the  
418 Targeted-Type scenario (Figure 4). Recent work on identifying when predictions will be  
419 extrapolation or interpolation suggests that this can be done by either examining distributions of  
420 predictor variables or comparing predictive variance at out-of-sample locations to a threshold  
421 (e.g. maximum predictive variance) based on in-sample locations (Bartley et al. In Press, Conn et  
422 al. 2015). As our example shows, not all targeted sampling designs will result in extrapolation  
423 and it may be acceptable to include data from such targeted efforts in larger, compiled datasets.

424 We also anticipated that stratified random sampling would result in better predictive  
425 performance than random sampling. There is intuitive appeal to stratified random sampling  
426 designs, particularly given the large amounts of ecological variation that exist at the regional  
427 scale (e.g., Cheruvilil et al. 2013, Lapierre et al. 2018). However, we did not find this to be the  
428 case, perhaps because LAGOS-NE includes a large sample size of 6,784 lakes (Table 3) that are  
429 relatively evenly distributed such that they capture the large geographic gradients that are present  
430 across the study area. It is also important to recognize that stratified random sampling design is  
431 better than random sampling design only if the strata used are ecologically relevant. For  
432 example, some sampling designs stratify by ecosystem size or area (e.g., U.S. EPA 2017).  
433 However, we did not include lake area as a stratum because it is not related to lake characteristics  
434 in LAGOS-NE (Stanley et al. 2019). Although we did not find stratified random designs to  
435 improve predictions over simple random designs, there may be other ways to stratify lakes that  
436 we did not consider here; further, there may be other ecosystem types, locations, or uses of  
437 macroscale monitoring data that require stratification.

438 Finally, we expected that lake nutrients would be better predicted than productivity  
439 because nutrients are more directly related to landscape context characteristics (Wagner and  
440 Schliep 2018) and exhibit a stronger spatial structure than CHL in our study area (Lapierre et al.  
441 2018). This expectation was supported by lower  $R^2$ s and higher RMSEs for CHL than for the  
442 nutrients. In fact, the errors may be enough to suggest an alternate trophic state (e.g., the

443 predicted value could be beyond a trophic threshold between mesotrophic and eutrophic). These  
444 results suggest that there is no one best sampling design for all response variables and that  
445 multiple metrics should be used when evaluating model predictive performance. The different  
446 diagnostic metrics also suggest that there are some subtle differences depending on which metric  
447 is used, and caution should be made in interpreting the results when selecting a sampling design  
448 based on one model performance metric alone.

449 Our conclusions should be interpreted within the context of the data used to conduct the  
450 research, specifically regarding the type of database, the study area, and the sample sizes used.  
451 This research was conducted using a compiled database of 87 disparate lake water quality  
452 datasets, many of which were sampled by individual U.S. state agencies (LAGOS-NE; Soranno  
453 et al. 2015, 2017). Consequently, sample lakes in LAGOS-NE were selected using a variety of  
454 different sampling strategies. In particular, sampled lakes tend to be larger and more connected  
455 than all lakes in the study area (Stanley et al. 2019). Therefore, lakes with *in situ* measurements  
456 in LAGOS-NE may not completely represent all lakes within the study area and the mimicked  
457 random sampling designs are not truly random (i.e., random selection from ~4,000 to 8,000 lakes  
458 with lake nutrients and productivity data rather than random selection from ~51,000 lakes in the  
459 census population). However, we believe that the sampled lakes in LAGOS-NE can provide a  
460 good approximation of all lakes in this geographic extent for three important reasons: (1)  
461 Because LAGOS-NE contains sampled lakes that vary widely by lake type, region, and  
462 ecological contexts, it contains sufficient variation in predictors and responses to effectively  
463 build predictive models; (2) a prior resampling exercise that corrected for the surface area  
464 sampling bias that we know is present in LAGOS-NE did not substantially change the statistical  
465 distributions of total nutrients and productivity (Stanley et al. 2019); and, (3) the combined  
466 sample sizes are likely large enough that any existing biases due to individual program sampling  
467 designs would have only minor effect on model performance.

#### 468 **Implications for macroscale sampling designs**

469 Macroscale monitoring programs often use either a stratified random design or targeted  
470 sampling. Our results from LAGOS-NE, which includes a variety of sampling designs, suggest  
471 that predictions from targeted sampling designs may sometimes perform similarly to those from  
472 random sampling designs. Thus, there is potential to include these datasets that were created to  
473 answer particular questions or to address specific environmental problems in compiled datasets  
474 because the bias associated with these data have only minimal effect on prediction errors. This

475 fact will be especially true when data from targeted sampling designs make up a small proportion  
476 of the total compiled ecosystem data, resulting in differences between the distributions of  
477 training and test datasets (i.e., extrapolation). Moreover, because it is unlikely that a large  
478 number of datasets will have exactly the same sample biases, assembling multiple datasets  
479 should tend to minimize the impact of any one dataset collected for one particular reason on  
480 prediction. Therefore, the use of such secondary datasets compiled from multiple sources, as was  
481 done for LAGOS-NE, is useful for macroscale prediction of ecosystem characteristics. Based on  
482 our results using lake nutrients and productivity, we make two specific suggestions for  
483 optimizing sampling designs at macroscales.

484 *To stratify or not?* Our results suggest that it may not be necessary to stratify when a  
485 relatively large sample size is feasible and relevant strata for prediction are either not present or  
486 unknown. Macroscale monitoring programs generally sample less than 20% of ecosystems, and  
487 sometimes as little as 1%. In comparison, LAGOS-NE includes 8 to 13% of all lakes  $\geq 4$  ha,  
488 depending on the response variable. And, for a stratified random design to be effective, the  
489 stratification must account for some variation in the ecosystem characteristics of interest such  
490 that resulting predictions are interpolations (predictions within model space) as opposed to  
491 extrapolation (predictions outside of model space). We tested two commonly used approaches  
492 for stratification that have been shown to capture variation in lake nutrients and productivity -  
493 regions (e.g., Cheruvilil et al. 2013) and local lake and watershed characteristics (e.g., Collins et  
494 al. 2018) - but were unable to document substantial improvements in model performance over  
495 simple random designs. This fact is likely because the relatively large number of lakes spread  
496 across wide environmental gradients in LAGOS-NE resulted in similar distributions of training  
497 and testing data (Figure 4) such that strata were not necessary to effectively capture variation in  
498 predictors and responses. Therefore, predictive performance was not substantially improved by  
499 adding stratification to a simple random sampling design.

500 Further, it is not likely that a single stratification design would adequately capture the  
501 complexity of all ecosystem characteristics, particularly when one considers biological, physical,  
502 and chemical characteristics in diverse ecosystems. Because the key characteristics that are most  
503 beneficial to use as strata will vary by response variable, it may be more effective to increase the  
504 total sample size across the study area rather than to spread samples across strata. Such relatively  
505 large and distributed sampling should help to increase predictive performance. The relative  
506 performance of simple random versus stratified random designs warrants testing in other



507 settings, for other macroscale datasets, and for other ecosystem characteristics to test the  
508 generality of our results. For example, the need for stratification may become more important as  
509 landscapes become more heterogeneous or vary across strata and as sample sizes drastically  
510 decrease, resulting in sampled ecosystems being less likely to represent a large proportion of the  
511 total landscapes or ecosystems within a study area.

512 *Space or time?* Our study examined the macroscale spatial predictions of lake nutrients  
513 and productivity by leveraging the broad spatial gradients in the LAGOS-NE database. In fact,  
514 an analysis of LAGOS-NE data using annual time scales across several decades found that  
515 spatial variation of lake nutrients and productivity far exceeded temporal variation (Soranno et  
516 al. 2019). However, if the goal is to predict responses of all ecosystems across regions and  
517 continents to a range of global change stressors, then making predictions across both space and  
518 time is essential (Janousek et al. 2019). Unfortunately, there are few spatially- and temporally-  
519 extensive datasets and the costs of long-term monitoring through both time and space are  
520 untenable. Thus, for new macroscale sampling programs, we recommend first capturing the  
521 existing spatial variation in predictor and response variables by sampling across the full range of  
522 ecological contexts present across a study area. Then, once sufficient spatial variation is  
523 captured, resources could be directed towards a smaller number of systems that are repeatedly  
524 sampled to capture temporal variation. By combining the use of secondary datasets that have  
525 excellent spatial coverage across a range of ecological context settings with sampling designs  
526 focused on filling in gaps in the temporal domain, macroscale studies will be able to inform a  
527 wide range of science questions and policy goals related to forecasting the effects of global  
528 change on ecosystem characteristics.

529

## 530 **ACKNOWLEDGMENTS**

531 Author contributions are as follows. PAS and KSC are co-lead authors and contributed equally to  
532 the manuscript by leading the conceptualization and writing of the manuscript. After the co-  
533 leads, there are 4 groups of authors in decreasing level of contribution, with authors listed in  
534 alphabetical order within each group. 1) QW and BL performed the analysis, with 2) P-NT and  
535 JZ as supervisors. 3) KBSK, IMM, and JS performed database queries, summaries, created tables  
536 and figures, and the code repository. 4) The remaining authors, in addition to those in groups 1-3,  
537 contributed to the development, editing, and writing of the paper. The authors declare that they  
538 have no conflict of interest. Further, we wish to thank Autumn Poisson and all participants from

539 the 2018 Continental Limnology Project Workshop at Pennsylvania State University, including  
540 Emily Stanley, Nicole Smith, Nathan Wikle, Sarah Collins, and Claire Boudreau. Thanks to  
541 Meredith and Justin Holgerson for their feedback and providing information about the FIA  
542 program, to Allie Shoffner for her editorial suggestions, and to the anonymous reviewers whose  
543 suggestions improved this manuscript. Funding was provided by the US NSF Macrosystems  
544 Biology Program grants, DEB-1638679; DEB-1638550, DEB-1638539, DEB-1638554. PAS  
545 was also supported by USDA National Institute of Food and Agriculture Hatch Project, Grant  
546 Number: 176820. Any use of trade, firm, or product names is for descriptive purposes only and  
547 does not imply endorsement by the U.S. Government.

548

#### 549 **LITERATURE CITED**

- 550 Bartley, M. L., E. M. Hanks, E. M. Schliep, P. A. Soranno, and T. Wagner. In Press. Identifying  
551 and characterizing extrapolation in multivariate response data. PLoS ONE.
- 552 Breiman, L. 2001. Random Forests. *Machine Learning* 45:5–32.
- 553 Cheruvilil, K. S., P. A. Soranno, K. E. Webster, and M. T. Bremigan. 2013. Multi-scaled drivers  
554 of ecosystem state: quantifying the importance of the regional spatial scale. *Ecological*  
555 *Applications* 23:1603–1618.
- 556 Collins, S. M., S. K. Oliver, J. Lapierre, E. H. Stanley, J. R. Jones, T. Wagner, and P. A.  
557 Soranno. 2017. Lake nutrient stoichiometry is less predictable than nutrient  
558 concentrations at regional and sub-continental scales. *Ecological Applications* 27:1529–  
559 1540.
- 560 Conn, P. B., D. S. Johnson, and P. L. Boveng. 2015. On Extrapolating Past the Range of  
561 Observed Data When Making Statistical Predictions in Ecology. PLoS ONE  
562 10:e0141416.
- 563 Dietze, M. C., A. Fox, L. M. Beck-Johnson, J. L. Betancourt, M. B. Hooten, C. S. Jarnevich, T.  
564 H. Keitt, M. A. Kenney, C. M. Laney, L. G. Larsen, H. W. Loescher, C. K. Lunch, B. C.  
565 Pijanowski, J. T. Randerson, E. K. Read, A. T. Tredennick, R. Vargas, K. C. Weathers,  
566 and E. P. White. 2018. Iterative near-term ecological forecasting: Needs, opportunities,  
567 and challenges. *Proceedings of the National Academy of Sciences* 115:1424–1432.
- 568 Gaiji, S., V. Chavan, A. H. Ariño, J. Otegui, D. Hobern, R. Sood, and E. Robles. 2013. Content  
569 assessment of the primary biodiversity data published through GBIF network: Status,  
570 challenges and potentials. *Biodiversity Informatics* 8.

571 Heffernan, J. B., P. A. Soranno, M. J. Angilletta, L. B. Buckley, D. S. Gruner, T. H. Keitt, J. R.  
572 Kellner, J. S. Kominoski, A. V. Rocha, J. Xiao, T. K. Harms, S. J. Goring, L. E. Koenig,  
573 W. H. McDowell, H. Powell, A. D. Richardson, C. A. Stow, R. Vargas, and K. C.  
574 Weathers. 2014. Macrosystems ecology: understanding ecological patterns and processes  
575 at continental scales. *Frontiers in Ecology and the Environment* 12:5–14.

576 Janousek, W. M., B. A. Hahn, and V. J. Dreitz. 2019. Disentangling monitoring programs:  
577 design, analysis, and application considerations. *Ecological Applications* 0:e01922.

578 Lapiere, J.-F., S. M. Collins, D. A. Seekell, K. S. Cheruvilil, P.-N. Tan, N. K. Skaff, Z. E.  
579 Taranu, C. E. Fergus, and P. A. Soranno. 2018. Similarity in spatial structure constrains  
580 ecosystem relationships: Building a macroscale understanding of lakes. *Global Ecology  
581 and Biogeography* 27:1251–1263.

582 Liaw, A., and M. Wiener. 2002. Classification and regression by randomForest. *R News* 2:5.

583 Lohr, S. L. 2019. *Sampling: Design and Analysis*. 2 edition. Routledge. New York, NY.

584 Miller, J. R., M. G. Turner, E. A. H. Smithwick, C. L. Dent, and E. H. Stanley. 2004. Spatial  
585 extrapolation: The science of predicting ecological patterns and processes. *BioScience*  
586 54:310–320.

587 Olsen, A. R., J. Sedransk, D. Edwards, C. A. Gotway, W. Liggett, S. Rathbun, K. H. Reckhow,  
588 and L. J. Young. 1999. Statistical Issues for Monitoring Ecological and Natural  
589 Resources in the United States. *Environmental Monitoring and Assessment* 54:1–45.

590 Pedregosa, F., G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P.  
591 Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, and D. Cournapeau. 2011.  
592 Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*  
593 12:2825–2830.

594 Peters, D. P. C., N. D. Burruss, L. L. Rodriguez, D. S. McVey, E. H. Elias, A. M. Pelzel-  
595 McCluskey, J. D. Derner, T. S. Schrader, J. Yao, S. J. Pauszek, J. Lombard, S. R. Archer,  
596 B. T. Bestelmeyer, D. M. Browning, C. W. Brungard, J. L. Hatfield, N. P. Hanan, J. E.  
597 Herrick, G. S. Okin, O. E. Sala, H. Savoy, and E. R. Vivoni. 2018. An integrated view of  
598 complex landscapes: A big data-model integration approach to transdisciplinary science.  
599 *BioScience* 68:653–669.

600 Poisson, A. C., I. M. McCullough, K. S. Cheruvilil, K. C. Elliott, J. A. Latimore, and P. A.  
601 Soranno. 2019. Quantifying the contribution of citizen science to broad-scale ecological  
602 databases. *Frontiers in Ecology and the Environment*. <https://doi.org/10.1002/fee.2128>

603 Rask, M., M. Olin, and J. Ruuhijärvi. 2010. Fish-based assessment of ecological status of Finnish  
604 lakes loaded by diffuse nutrient pollution from agriculture. *Fisheries Management and*  
605 *Ecology* 17:126–133.

606 Read, E. K., V. P. Patil, S. K. Oliver, A. L. Hetherington, J. A. Brentrup, J. A. Zwart, K. M.  
607 Winters, J. R. Corman, E. R. Nodine, R. I. Woolway, H. A. Dugan, A. Jaimes, A. B.  
608 Santoso, G. S. Hong, L. A. Winslow, P. C. Hanson, and K. C. Weathers. 2015. The  
609 importance of lake-specific characteristics for water quality across the continental United  
610 States. *Ecological Applications* 25:943–955.

611 Sauer, J. R., W. A. Link, J. E. Fallon, K. L. Pardieck, and D. J. Ziolkowski. 2013. The North  
612 American Breeding Bird Survey 1966–2011: Summary analysis and species accounts.  
613 *North American Fauna*:1–32.

614 Seaber, P. R., F. P. Kapinos, and G. L. Knapp. 1987. Hydrologic unit maps. USGS Numbered  
615 Series, U.S. G.P.O.,.

616 Soranno, P. A., L. C. Bacon, M. Beauchene, K. E. Bednar, E. G. Bissell, C. K. Boudreau, M. G.  
617 Boyer, M. T. Bremigan, S. R. Carpenter, J. W. Carr, K. S. Cheruvilil, S. T. Christel, M.  
618 Claucherty, S. M. Collins, J. D. Conroy, J. A. Downing, J. Dukett, C. E. Fergus, C. T.  
619 Filstrup, C. Funk, M. J. Gonzalez, L. T. Green, C. Gries, J. D. Halfman, S. K. Hamilton,  
620 P. C. Hanson, E. N. Henry, E. M. Herron, C. Hockings, J. R. Jackson, K. Jacobson-  
621 Hedin, L. L. Janus, W. W. Jones, J. R. Jones, C. M. Keson, K. B. S. King, S. A.  
622 Kishbaugh, J.-F. Lapierre, B. Lathrop, J. A. Latimore, Y. Lee, N. R. Lottig, J. A. Lynch,  
623 L. J. Matthews, W. H. McDowell, K. E. B. Moore, B. P. Neff, S. J. Nelson, S. K. Oliver,  
624 M. L. Pace, D. C. Pierson, A. C. Poisson, A. I. Pollard, D. M. Post, P. O. Reyes, D. O.  
625 Rosenberry, K. M. Roy, L. G. Rudstam, O. Sarnelle, N. J. Schuldt, C. E. Scott, N. K.  
626 Skaff, N. J. Smith, N. R. Spinelli, J. J. Stachelek, E. H. Stanley, J. L. Stoddard, S. B.  
627 Stopyak, C. A. Stow, J. M. Tallant, P.-N. Tan, A. P. Thorpe, M. J. Vanni, T. Wagner, G.  
628 Watkins, K. C. Weathers, K. E. Webster, J. D. White, M. K. Wilmes, and S. Yuan. 2017.  
629 LAGOS-NE: a multi-scaled geospatial and temporal database of lake ecological context  
630 and water quality for thousands of US lakes. *GigaScience* 6:1–22.

631 Soranno, P. A., E. G. Bissell, K. S. Cheruvilil, S. T. Christel, S. M. Collins, C. E. Fergus, C. T.  
632 Filstrup, J.-F. Lapierre, N. R. Lottig, S. K. Oliver, C. E. Scott, N. J. Smith, S. Stopyak, S.  
633 Yuan, M. T. Bremigan, J. A. Downing, C. Gries, E. N. Henry, N. K. Skaff, E. H. Stanley,  
634 C. A. Stow, P.-N. Tan, T. Wagner, and K. E. Webster. 2015. Building a multi-scaled

635 geospatial temporal ecology database from disparate data sources: fostering open science  
636 and data reuse. *GigaScience* 4:28.

637 Soranno, P. A., and K. S. Cheruvilil. 2017a. LAGOS-NE-GEO v1.05: A module for LAGOS-  
638 NE, a multi-scaled geospatial and temporal database of lake ecological context and water  
639 quality for thousands of U.S. Lakes: 1925-2013.

640 Soranno, P. A., and K. S. Cheruvilil. 2017b. LAGOS-NE-LIMNO v1.087.1: A module for  
641 LAGOS-NE, a multi-scaled geospatial and temporal database of lake ecological context  
642 and water quality for thousands of U.S. Lakes: 1925-2013. DOI:  
643 10.6073/pasta/b1b93ccf3354a7471b93ecca484d506.

644 Soranno, P. A., T. Wagner, S. M. Collins, J.-F. Lapierre, N. R. Lottig, and S. K. Oliver. 2019.  
645 Spatial and temporal variation of ecosystem properties at macroscales. *Ecology Letters*  
646 22:1587–1598.

647 Stanley, E. H., S. M. Collins, N. R. Lottig, S. K. Oliver, K. E. Webster, K. S. Cheruvilil, and P.  
648 A. Soranno. 2019. Biases in lake water quality sampling and implications for macroscale  
649 research. *Limnology and Oceanography* 64:1572–1585.

650 Stone, M. 1974. Cross-validatory choice and assessment of statistical predictions. *Journal of the*  
651 *Royal Statistical Society. Series B (Methodological)* 36:111–147.

652 Swetnam, T. W. 1993. Fire history and climate change in giant Sequoia groves. *Science*  
653 262:885–889.

654 Thompson, S. K. 2012. *Sampling*. 3 edition. Wiley, Hoboken, N.J.

655 Urquhart, N. S., S. G. Paulsen, and D. P. Larsen. 1998. Monitoring for policy-relevant regional  
656 trends over time. *Ecological Applications* 8:246–257.

657 US Environmental Protection Agency. 1975. *National Eutrophication Survey Methods 1973–*  
658 *1976 (Working Paper No. 175)*, Tech. rep. United States Environmental Protection  
659 Agency, Office of Research and Development, Corvallis, OR, USA.

660 U.S. Environmental Protection Agency Office of Wetlands, Oceans and Watersheds Office of  
661 Research and Development. 2017. *National Lakes Assessment 2012: Technical Report*.  
662 U.S. Environmental Protection Agency, Washington DC.

663 Ver Hoef, J. M. 2008. Spatial methods for plot-based sampling of wildlife populations.  
664 *Environmental and Ecological Statistics* 15:3–13.

665 Wagner, T., N. R. Lottig, E. M. Schliep, J. J. Stachelek, K. S. Cheruvilil, and S. M. Collins.  
666 2020. Increasing accuracy of nutrient predictions in thousands of lakes by leveraging

667 water clarity data by citizen scientists. *Limnology and Oceanography Letters*.  
 668 <https://aslopubs.onlinelibrary.wiley.com/doi/full/10.1002/lol2.10134>

669 Wagner, T., and E. M. Schliep. 2018. Combining nutrient, productivity, and landscape-based  
 670 regressions improves predictions of lake nutrients and provides insight into nutrient  
 671 coupling at macroscales. *Limnology and Oceanography* 63:2372–2383.

672 Ward, J. H. 1963. Hierarchical grouping to optimize an objective function. *Journal of the*  
 673 *American Statistical Association* 58:236–244.

674 Webb, E. L., D. A. Friess, K. W. Krauss, D. R. Cahoon, G. R. Guntenspergen, and J. Phelps.  
 675 2013. A global standard for monitoring coastal wetland vulnerability to accelerated sea-  
 676 level rise. *Nature Climate Change* 3:458–465.

677 Zhao, S., N. Pederson, L. D’Orangeville, J. HilleRisLambers, E. Boose, C. Penone, B. Bauer, Y.  
 678 Jiang, and R. D. Manzanedo. 2019. The International Tree-Ring Data Bank (ITRDB)  
 679 revisited: Data availability and global ecological representativity. *Journal of*  
 680 *Biogeography* 46:355–368.

681 Zhou, Z.-H. 2012. Ensemble methods: Foundations and algorithms. Chapman & Hall Inc, Boca  
 682 Raton, FL.

683  
 684  
 685 **DATA AVAILABILITY**

686  
 687 Data and code are available in Zenodo at: <http://doi.org/10.5281/zenodo.3606832>

688  
 689 **TABLES**

690  
 691 **Table 1:** Summary of response variables (minimum, maximum, median, mean, 25<sup>th</sup> and 75<sup>th</sup>  
 692 percentiles.

Response Variable	Units	n	Min	25th	Median	Mean	75th	Max
Total Phosphorus	µg/l	5896	0	10	16	39.9	34	1184
Total Nitrogen	µg/l	4253	0	380	600	944.3	990	20574
Chlorophyll-a	µg/l	6784	0	2.6	5	16.28	13	553.4

693  
 694

695  
 696  
 697  
 698  
 699  
 700  
 701  
 702

**Table 2:** Summary of the predictor variables (minimum, maximum, median, mean, 25<sup>th</sup> and 75<sup>th</sup> percentiles) at the local (lake and watershed) and regional scales. Note, lake connectivity is a categorical variable with 3 categories (Isolated, DR\_Stream, and DR\_LakeStream). Lake and watershed complexity refer to lake and watershed boundary complexity factor, respectively, which are measures of reticulation, N dep refers to nitrogen deposition in 1990, and N dep difference refers to N deposition in 1990 minus that in 2010.

Predictor Variable	Units	Min	25th	Median	Mean	75th	Max
<b>Local</b>							
lake connectivity <sup>1</sup>	NA	NA	NA	NA	NA	NA	NA
lake water clarity <sup>1</sup>	m	0	1.30	2.40	2.75	3.80	18.25
lake max depth <sup>1</sup>	m	0.30	4.60	8.53	10.84	14.02	198.4
lake complexity <sup>1</sup>	NA	1.00	1.40	1.75	2.11	2.35	30.27
lake elevation <sup>2</sup>	m	0	241.1	323.9	316.5	412.1	1038.6
watershed wetland <sup>3</sup>	%	0.00	2.42	7.23	12.28	17.77	93.08
watershed complexity <sup>1</sup>	NA	1.21	2.02	2.37	2.59	2.85	25.49
watershed lake ratio <sup>1</sup>	NA	0.01	3.88	8.31	42.63	21.03	53517.4
watershed stream density <sup>1</sup>	m/ha	0	0	3.08	4.54	7.57	71.77
watershed forest <sup>3</sup>	%	0	23.70	53.8	49.91	75.05	100
watershed road density <sup>4</sup>	m/ha	0	14.50	24.35	30.96	39.36	262.66
<b>Regional</b>							
baseflow mean <sup>5</sup>	%	14.18	47.92	52.62	52.08	58.44	78.83
runoff mean <sup>5</sup>	in/year	2.80	7.26	10.65	13.21	22.59	26.95
agriculture <sup>3</sup>	%	1.79	5.67	26.33	28.64	34.07	78.66
temperature mean <sup>6</sup>	°C	3.46	5.44	6.15	6.83	8.17	15.40
precipitation mean <sup>6</sup>	mm	606.60	714	839.3	910.3	1106.8	1282.7
N dep mean <sup>7</sup>	kg/ha	2.68	4.37	5.36	5.27	5.99	8.67
N dep difference <sup>7</sup>	kg/ha	-1.49	-0.11	1.47	1.30	2.47	4.66

703  
 704  
 705  
 706  
 707

<sup>1</sup> National Hydrography Dataset (NHD) 2013; and Soranno et al. 2015  
<sup>2</sup> USGS National Elevation Dataset (NED); 2013  
<sup>3</sup> National Land Cover Database (NLCD); 2006  
<sup>4</sup> United States Census TIGER roads data; 2013

708 <sup>5</sup>United States Geological Survey (USGS); 1951-1980  
 709 <sup>6</sup>PRISM climate group 30-year normal; 1981-2010  
 710 <sup>7</sup>National Atmospheric Deposition Program; 1990-2010

711  
 712 **Table 3.** The number of lakes in the training and testing datasets for each of the seven sampling  
 713 scenarios and three response variables. For all scenarios except (g), these are average numbers of  
 714 lakes over multiple subsets of training and testing data. The numbers in parentheses are the  
 715 percent of the total lake population  $\geq 4$  ha comprised for each scenario and response variable  
 716 combination.

Sampling Scenario	TP		TN		CHL	
	Training	Testing	Training	Testing	Training	Testing
(a) Random-Large	4,422 (9 %)	1,474 (3 %)	3,190 (6 %)	1,063 (2 %)	5,088 (10 %)	1,696 (3 %)
(b) Random-Small	1,474 (3 %)	4,422 (9 %)	1,063 (2 %)	3,190 (6 %)	1,696 (3 %)	5,088 (10 %)
(c) Stratified-Type	1,474 (3 %)	4,422 (9 %)	1,063 (2 %)	3,190 (6 %)	1,696 (3 %)	5,088 (10 %)
(d) Stratified-Region	1,474 (3 %)	4,422 (9 %)	1,063 (2 %)	3,190 (6 %)	1,696 (3 %)	5,088 (10 %)
(e) Targeted-Type	2,869 (6 %)	3,024 (6 %)	2,079 (4 %)	2,171 (4 %)	3,393 (7 %)	3,379 (7 %)
(f) Targeted-Region	2,927 (6 %)	2,927 (6 %)	2,127 (4 %)	2,127 (4 %)	3,392 (7 %)	3,392 (7 %)
(g) Targeted-AgRegion	4,422 (9 %)	1,474 (3 %)	3,190 (6 %)	1,063 (2 %)	5,088 (10 %)	1,696 (3 %)

717  
 718  
 719  
 720 **FIGURE LEGENDS**  
 721  
 722 **Figure 1.** Conceptual figure depicting three types of sampling strategies (1-3) used for selecting  
 723 ecosystems to sample at macroscales. Underneath each type is a description of the assumptions  
 724 underlying the resulting models. In all seven depictions (a-g), there are ecosystems that are used



725 to build predictive models (training dataset; blue circles) and ecosystems that are used to test the  
726 predictive models (test, orange circles). From left to right: 1. Random sampling designs whereby  
727 ecosystems are chosen completely randomly from the sample extent; predictive models for  
728 unsampled ecosystems are assumed to be interpolation, if sample size is sufficient. 2. Stratified  
729 random sampling designs whereby ecosystems are first stratified by ecosystem type (top) or their  
730 location within ecological regions (regions depicted by dark lines, bottom) that are thought to  
731 drive variation among ecosystems and second, ecosystems are selected randomly within those  
732 strata; predictive models for unsampled ecosystems are assumed to be interpolation, if the strata  
733 are ecologically relevant and sample size is sufficient. 3. Targeted sampling whereby particular  
734 types of ecosystems (top), particular ecological regions (middle), or regions with particular land  
735 uses (bottom) are targeted for sampling in order to answer a particular question; predictive  
736 models for unsampled ecosystems are assumed to be extrapolation. Black lower-case letters  
737 relate to the seven scenarios used in this study that are described and depicted throughout.

738

739 **Figure 2.** Map of lakes color coded by water clarity measured as Secchi disk depth (m), colored  
740 by percentile. Gray lines delineate regions.

741

742 **Figure 3.** Boxplots showing the model predictive performance of each scenario indicated by  
743 letters (X axis labels, letters as per Fig. 1) as measured by predictive  $R^2$  (A), root mean square  
744 error (RMSE; B), and median relative absolute error (MRAE; C). The colors signify the different  
745 types of sampling strategies: random (yellow), stratified (green), and targeted (blue). Y-axis  
746 scales are truncated for better visualization.

747

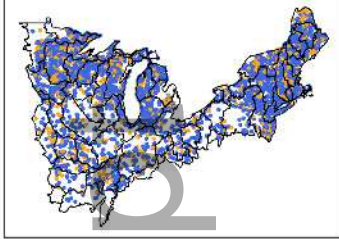
748 **Figure 4.** Density plots showing the distribution of data in the training (blue) and testing  
749 (orange) dataset for each sample design scenario (a-g as per Fig. 1) and for left to right: TP  
750 ( $\mu\text{g/L}$ ), TN ( $\mu\text{g/L}$ ), CHL ( $\mu\text{g/L}$ ), water clarity (m), lake maximum depth (m), and watershed  
751 percent forested. One randomly selected dataset for each sample design scenario is portrayed in  
752 this figure. The X-axis is truncated and axis labels are not shown to better visualize the majority  
753 of the data for best visual comparison of the training vs test datasets. The letters are as for Figure  
754 1.

# Author Manuscript

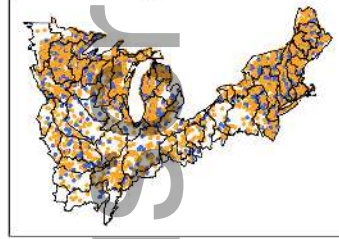
### 1. Random sampling

**Interpolation**, sampled ecosystems representative of unsampled ecosystems, if sample size is sufficient.

(a) Large training dataset



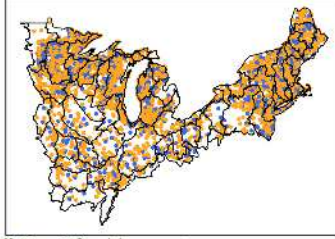
(b) Small training dataset



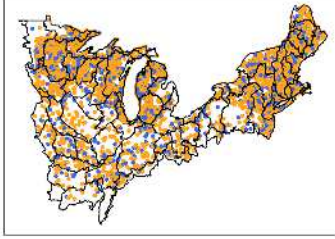
### 2. Stratified random sampling

**Interpolation**, sampled ecosystems representative of unsampled ecosystems, if strata ecologically relevant.

(c) Stratified by ecosystem type



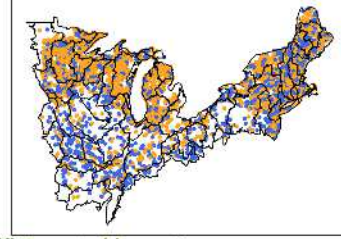
(d) Stratified by region



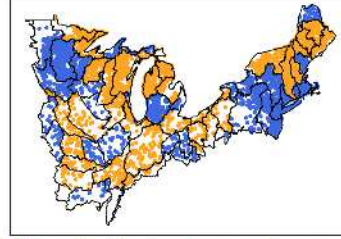
### 3. Targeted sampling

**Extrapolation**, sampled ecosystems not representative of unsampled ecosystems.

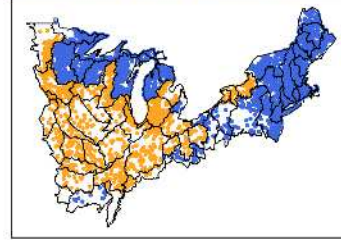
(e) Targeted by ecosystem type



(f) Targeted by region



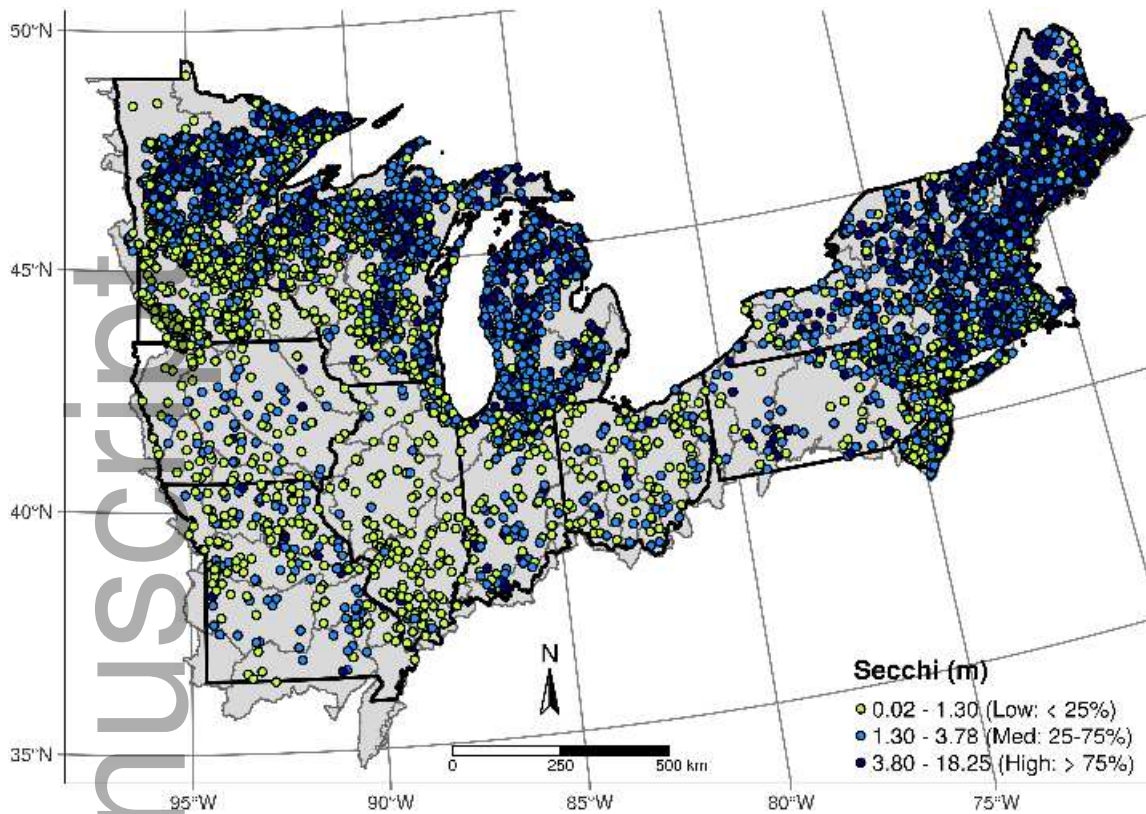
(g) Targeted by land use regions



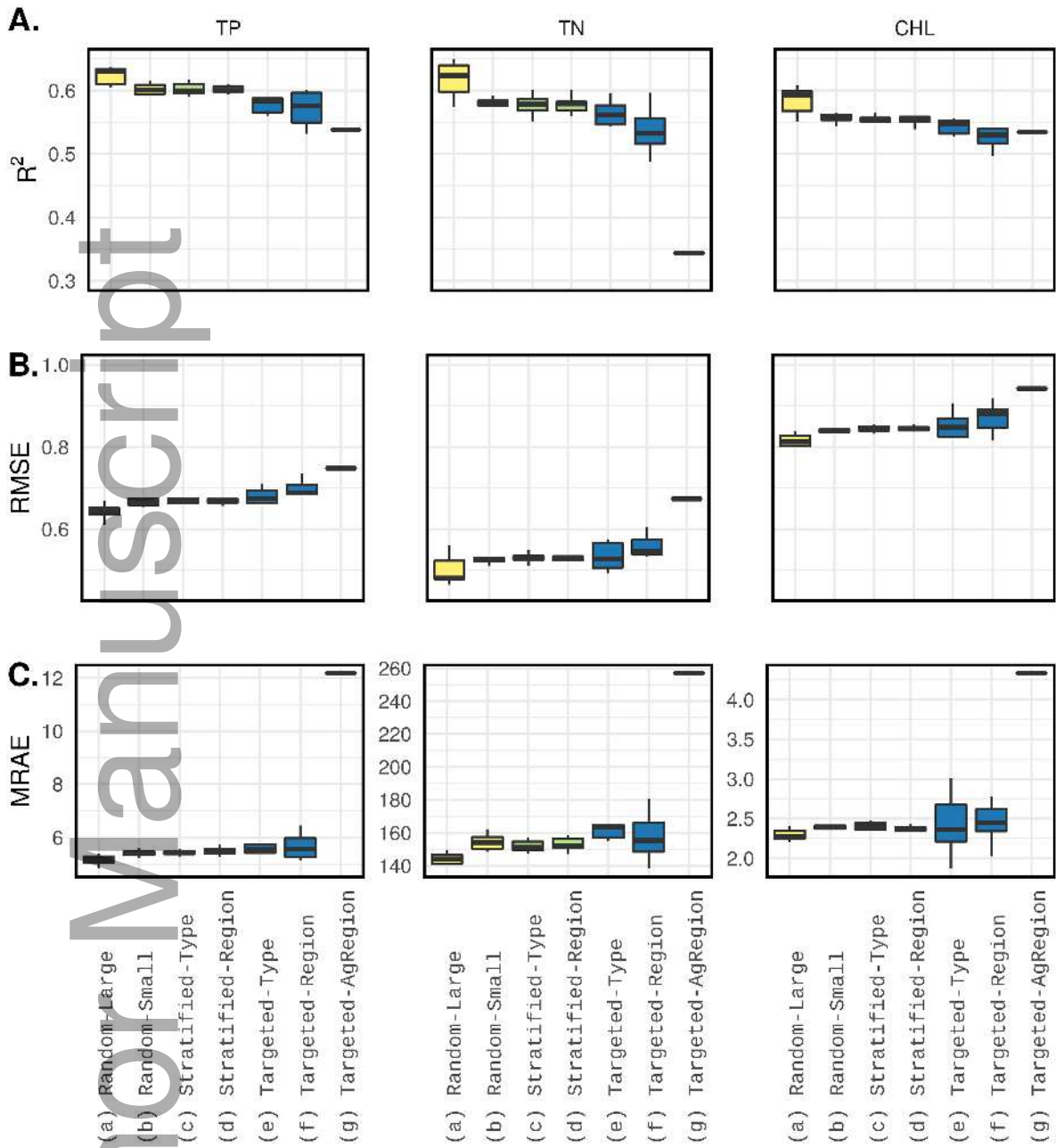
● Testing  
● Training

eap\_2123\_f1.tif

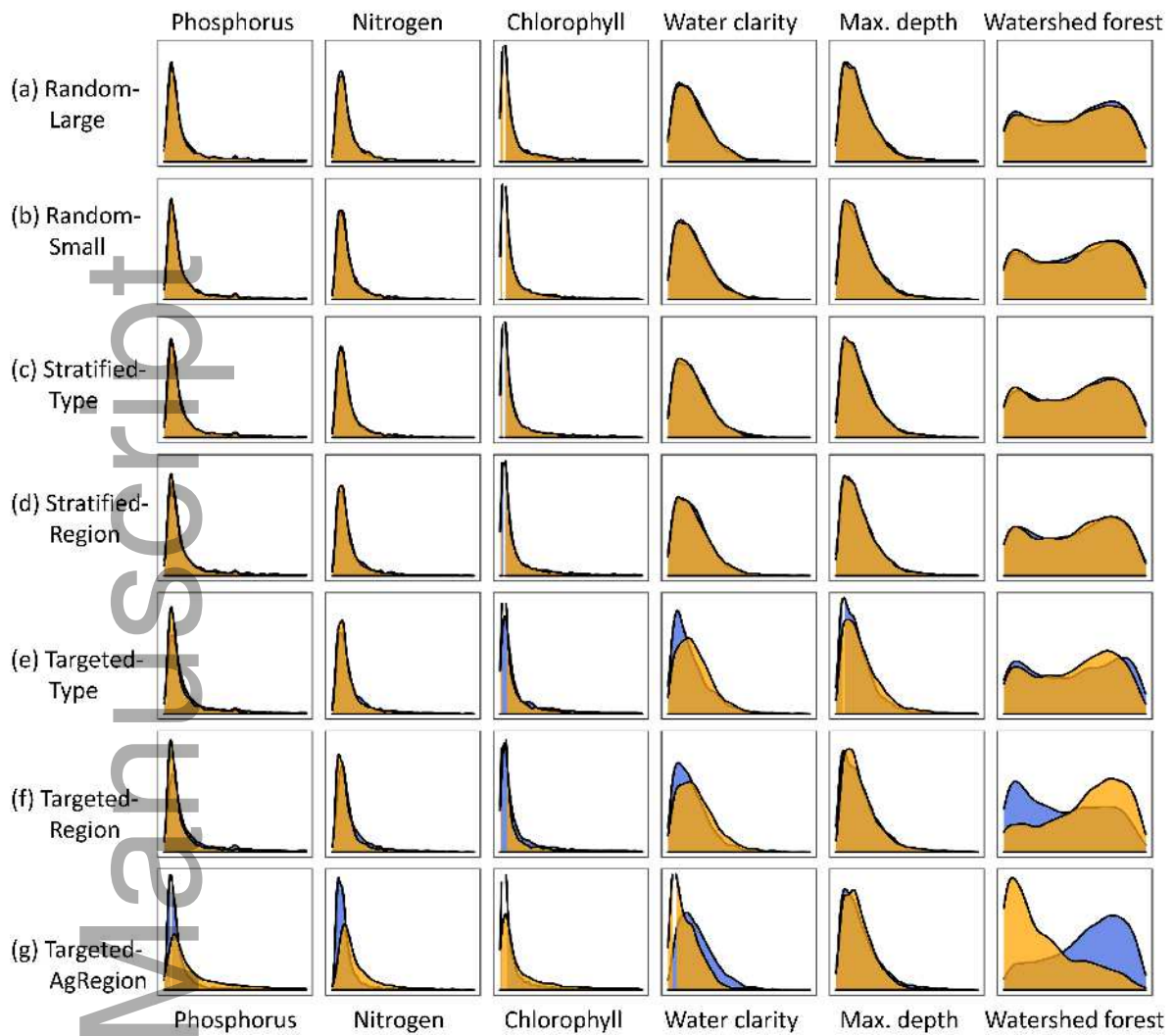
Author Manuscript



eap\_2123\_f2.tif



eap\_2123\_f3.tif



eap\_2123\_f4.tif

Author