

## DATA NOTE

## LAGOS-NE: a multi-scaled geospatial and temporal database of lake ecological context and water quality for thousands of US lakes

Patricia A. Soranno<sup>1,\*</sup>, Linda C. Bacon<sup>2</sup>, Michael Beauchene<sup>3</sup>, Karen E. Bednar<sup>4</sup>, Edward G. Bissell<sup>1</sup>, Claire K. Boudreau<sup>1</sup>, Marvin G. Boyer<sup>5</sup>, Mary T. Bremigan<sup>1</sup>, Stephen R. Carpenter<sup>6</sup>, Jamie W. Carr<sup>7</sup>, Kendra S. Cheruvilil<sup>1</sup>, Samuel T. Christel<sup>6</sup>, Matt Claucherty<sup>8</sup>, Sarah M. Collins<sup>6</sup>, Joseph D. Conroy<sup>9</sup>, John A. Downing<sup>10</sup>, Jed Dukett<sup>11</sup>, C. Emi Fergus<sup>12</sup>, Christopher T. Filstrup<sup>10</sup>, Clara Funk<sup>13</sup>, Maria J. Gonzalez<sup>14</sup>, Linda T. Green<sup>15</sup>, Corinna Gries<sup>6</sup>, John D. Halfman<sup>16</sup>, Stephen K. Hamilton<sup>17</sup>, Paul C. Hanson<sup>6</sup>, Emily N. Henry<sup>18</sup>, Elizabeth M. Herron<sup>19</sup>, Celeste Hockings<sup>20</sup>, James R. Jackson<sup>21</sup>, Kari Jacobson-Hedin<sup>22</sup>, Lorraine L. Janus<sup>23</sup>, William W. Jones<sup>24</sup>, John R. Jones<sup>25</sup>, Caroline M. Keson<sup>26</sup>, Katelyn B.S. King<sup>1</sup>, Scott A. Kishbaugh<sup>27</sup>, Jean-Francois Lapierre<sup>28</sup>, Barbara Lathrop<sup>29</sup>, Jo A. Latimore<sup>1</sup>, Yuehlin Lee<sup>30</sup>, Noah R. Lottig<sup>31</sup>, Jason A. Lynch<sup>13</sup>, Leslie J. Matthews<sup>32</sup>, William H. McDowell<sup>33</sup>, Karen E.B. Moore<sup>34</sup>, Brian P. Neff<sup>35</sup>, Sarah J. Nelson<sup>36</sup>, Samantha K. Oliver<sup>6</sup>, Michael L. Pace<sup>37</sup>, Donald C. Pierson<sup>38</sup>, Autumn C. Poisson<sup>1</sup>, Amina I. Pollard<sup>39</sup>, David M. Post<sup>40</sup>, Paul O. Reyes<sup>30</sup>, Donald O. Rosenberry<sup>41</sup>, Karen M. Roy<sup>42</sup>, Lars G. Rudstam<sup>43</sup>, Orlando Sarnelle<sup>1</sup>, Nancy J. Schuldt<sup>44</sup>, Caren E. Scott<sup>45</sup>, Nicholas K. Skaff<sup>1</sup>, Nicole J. Smith<sup>1</sup>, Nick R. Spinelli<sup>46</sup>, Joseph J. Stachelek<sup>1</sup>, Emily H. Stanley<sup>6</sup>, John L. Stoddard<sup>47</sup>, Scott B. Stopyak<sup>48</sup>, Craig A. Stow<sup>49</sup>, Jason M. Tallant<sup>50</sup>, Pang-Ning Tan<sup>51</sup>, Anthony P. Thorpe<sup>25</sup>,

Received: 2 June 2017; Revised: 8 September 2017; Accepted: 5 October 2017

© The Author 2017. Published by Oxford University Press. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

Michael J. Vanni<sup>52</sup>, Tyler Wagner<sup>53</sup>, Gretchen Watkins<sup>4</sup>,  
Kathleen C. Weathers<sup>54</sup>, Katherine E. Webster<sup>55</sup>, Jeffrey D. White<sup>56</sup>,  
Marcy K. Wilmes<sup>57</sup> and Shuai Yuan<sup>51</sup>

<sup>1</sup>Department of Fisheries and Wildlife, Michigan State University, East Lansing, MI 48824, USA, <sup>2</sup>Department of Environmental Protection, State of Maine, Augusta, ME 04330, USA, <sup>3</sup>Department of Energy and Environmental Protection, State of Connecticut, Hartford, CT 06106, USA, <sup>4</sup>Water Resources Program, Lac du Flambeau Tribal Natural Resources, Lac du Flambeau, WI, USA, <sup>5</sup>Environmental Planning, US Army Corps of Engineers, Kansas City, MO 64106, USA, <sup>6</sup>Center for Limnology, University of Wisconsin Madison, Madison, WI 53706 USA, <sup>7</sup>Office of Watershed Management, Massachusetts Department of Conservation and Recreation, West Boylston, MA 10583, USA, <sup>8</sup>Watershed Protection, Tipp of the Mitt Watershed Council, Petoskey, MI 49770, USA, <sup>9</sup>Division of Wildlife, Inland Fisheries Research Unit, Ohio Department of Natural Resources, Hebron, OH 43025, USA, <sup>10</sup>Large Lakes Observatory, University of Minnesota, Duluth, MN 55812 USA, <sup>11</sup>Adirondack Lake Survey Corporation, Ray Brook, NY 12977 USA, <sup>12</sup>National Research Council, US Environmental Protection Agency, Corvallis, OR 97333, USA, <sup>13</sup>Office of Air and Radiation, US Environmental Protection Agency, Washington, DC 20460, USA, <sup>14</sup>Department of Biology, Miami University, Oxford, OH 45056, USA, <sup>15</sup>Natural Resource Science, University of Rhode Island, Kingston, RI 02892 USA, <sup>16</sup>Geoscience, Hobart & William Smith Colleges, Geneva, NY 14456 USA, <sup>17</sup>Kellogg Biological Station, Michigan State University, Hickory Corners, MI 49060, USA, <sup>18</sup>Outreach and Engagement, Oregon State University, Corvallis, OR 97331, USA, <sup>19</sup>Watershed Watch, University of Rhode Island, Kingston, RI 02881, USA, <sup>20</sup>Natural Resource Department, Lac du Flambeau Band of Lake Superior Chippewa Indians, Lac du Flambeau, WI 54538, USA, <sup>21</sup>Department of Natural Resources, Cornell University, Bridgeport, NY, USA, <sup>22</sup>Office of Water Protection, Fond du Lac Reservation, Cloquet, MN 55720 USA, <sup>23</sup>Bureau of Water Supply, New York City Department of Environmental Protection, Valhalla, NY 10560, USA, <sup>24</sup>School of Public and Environmental Affairs, Indiana University, Bloomington, IN 47408, USA, <sup>25</sup>School of Natural Resources, University of Missouri, Columbia, MO, USA, <sup>26</sup>Natural Resource Department, Little Traverse Bay Bands of Odawa Indians, Harbor Springs, MI 49740, USA, <sup>27</sup>Division of Water, New York State Department of Environmental Conservation, Albany, NY 12233, USA, <sup>28</sup>Department of Biological Science, University of Montreal, Montreal Quebec, Canada, H3C 3J7, <sup>29</sup>Pennsylvania Department of Environmental Protection, State of Pennsylvania, Harrisburg, PA 17101 USA, <sup>30</sup>Office of Watershed Management, Massachusetts Department of Conservation and Recreation, Belchertown, MA 01007, USA, <sup>31</sup>Trout Lake Research Station, University of Wisconsin, Boulder Junction, WI 54512, USA, <sup>32</sup>Lakes and Ponds Program, Vermont Department of Environmental Conservation, Montpelier, VT 05620, USA, <sup>33</sup>Natural Resources and the Environment, University of New Hampshire, Durham, NH 03824, USA, <sup>34</sup>Water Quality Science and Research, New York City Department of Environmental Protection, Kingston, NY 12401, USA, <sup>35</sup>National Research Program, USGS, Denver CO 80225, USA, <sup>36</sup>School of Forest Resources, University of Maine, Orono, ME, USA, <sup>37</sup>Department of Environmental Science, University of Virginia, Charlottesville, VA 22904, USA, <sup>38</sup>Department of Ecology and Genetics, Uppsala University, Uppsala, Sweden, <sup>39</sup>Office of Water, US EPA, Washington, DC 20460, USA, <sup>40</sup>Ecology and Evolutionary Biology, Yale University, Connecticut 06511, USA, <sup>41</sup>National Research Program, USGS, Denver, CO 80225, USA, <sup>42</sup>Division of Air Resources, New York State Department of Environmental Conservation, Ray Brook, NY 12977, USA, <sup>43</sup>Department of Natural Resources, Cornell University, Ithaca, NY 14850, USA, <sup>44</sup>Environmental Program, Fond du Lac Band of Lake Superior Chippewa Indians, Cloquet, MN 55720, USA, <sup>45</sup>Aquatic Science, NEON, Boulder, CO 80301, USA, <sup>46</sup>Watershed Management, Lake Wallenpaupack Watershed Management District, Hawley, PA, USA, <sup>47</sup>Western Ecology Division, Office of Research and Development, US EPA, Corvallis, OR 97333, USA, <sup>48</sup>Technology Services, Eaton County, Charlotte, MI, USA, <sup>49</sup>Great Lakes Environmental Research Lab, NOAA, Ann Arbor, MI 47176, USA, <sup>50</sup>Biological Station, University of Michigan, Pellston, MI 49769, USA, <sup>51</sup>Computer Science and Engineering,

Michigan State University, East Lansing, MI 48824, USA, <sup>52</sup>Department of Zoology, Miami University, Oxford, OH 45056 USA, <sup>53</sup>Pennsylvania Cooperative Fish and Wildlife Research Unit, USGS, 402 Forest Resources Building, University Park, PA 16802, USA, <sup>54</sup>Cary Institute of Ecosystem Studies, Millbrook, NY, USA, <sup>55</sup>School of Natural Sciences, Trinity College, Dublin, Ireland, <sup>56</sup>Biology Department, Framingham State University, Framingham, MA 01702, USA and <sup>57</sup>Department of Environmental Quality, State of Michigan, Lansing, MI 48909, USA

\*Correspondence address. Patricia A. Soranno, PhD, Department of Fisheries and Wildlife, Michigan State University, East Lansing, MI 48824, USA. Tel: 517-884-1769; Fax: 517-432-1699; E-mail: [soranno@msu.edu](mailto:soranno@msu.edu)

## Abstract

Understanding the factors that affect water quality and the ecological services provided by freshwater ecosystems is an urgent global environmental issue. Predicting how water quality will respond to global changes not only requires water quality data, but also information about the ecological context of individual water bodies across broad spatial extents. Because lake water quality is usually sampled in limited geographic regions, often for limited time periods, assessing the environmental controls of water quality requires compilation of many data sets across broad regions and across time into an integrated database. LAGOS-NE accomplishes this goal for lakes in the northeastern-most 17 US states. LAGOS-NE contains data for 51 101 lakes and reservoirs larger than 4 ha in 17 lake-rich US states. The database includes 3 data modules for: lake location and physical characteristics for all lakes; ecological context (i.e., the land use, geologic, climatic, and hydrologic setting of lakes) for all lakes; and in situ measurements of lake water quality for a subset of the lakes from the past 3 decades for approximately 2600–12 000 lakes depending on the variable. The database contains approximately 150 000 measures of total phosphorus, 200 000 measures of chlorophyll, and 900 000 measures of Secchi depth. The water quality data were compiled from 87 lake water quality data sets from federal, state, tribal, and non-profit agencies, university researchers, and citizen scientists. This database is one of the largest and most comprehensive databases of its type because it includes both in situ measurements and ecological context data. Because ecological context can be used to study a variety of other questions about lakes, streams, and wetlands, this database can also be used as the foundation for other studies of freshwaters at broad spatial and ecological scales.

**Keywords:** lake eutrophication; nutrients; water quality; lake trophic state; ecological context; LAGOS-NE; open science; lake database

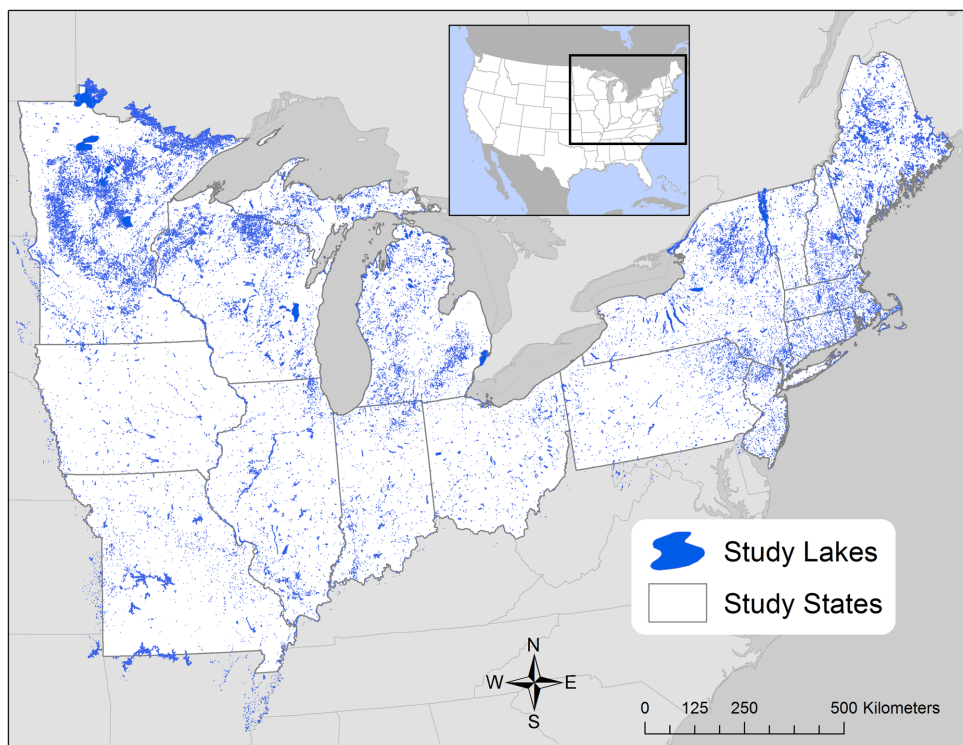
## Data Description

A major concern for water quality in freshwaters globally is cultural eutrophication, or excess nutrient inputs from human activities that lead to increased plant and algal growth. In many parts of the world, runoff from land, or nonpoint-source pollution, has replaced discharges of sewage, or point-source pollution, as the primary driver of lake and reservoir eutrophication [1]. In lakes and reservoirs, eutrophication is expected to become more widespread in the coming decades as the human population increases and climate and land use change commensurately, placing increasing pressures on freshwaters [2–4], although there is also recognition that eutrophication or its response to management actions does not progress in the same way in all lakes (e.g., [5–7]). Most research to understand lake nutrients and their effects on algae, plants, and aquatic food webs has been conducted in individual or small groups of lakes by studying the complex within-lake mechanisms that control responses to nutrients (e.g., [8, 9]). Such relationships and interactions have also been found to be influenced by the ecological context of lakes (i.e., the land use, geologic, climatic, and hydrologic setting of lakes), which varies by lake and region and is multi-scaled. In fact, it is not always clear whether local or regional ecological context matters more for predicting lake eutrophication (e.g., [10–12]). Therefore, determining the current extent of lake eutrophication and predicting how eutrophication will respond to future global change requires water quality data (e.g., nutrients, water clarity, and chlorophyll concen-

trations) and measures of lake ecological context across regions, the continent, and the globe (e.g., 13–15).

In practice, measures of water quality are often collected from a relatively small number of lakes within individual regions. In the United States, large investments have been made in water quality monitoring by federal, state, local, and tribal governments; and many, but not all, of the data sets have been placed in government data repositories such as the USGS National Water Information System (NWIS) and the USEPA Storage and Retrieval (STORET) database. Unfortunately, these data repositories do not currently allow us to study lake water quality at broad scales. Despite the large number of water quality records in these systems, a recent analysis of stream nutrient data obtained from NWIS, STORET, and more than 400 other organizations determined that more than half of the data records lacked the most critical metadata necessary to make the data usable (e.g., chemical form, parameter name, units) [16], and we would expect a similar result with lake data because they are typically treated similarly to stream nutrient data. In addition, STORET and NWIS do not include any measures of lake ecological context. Therefore, to study the controls of eutrophication specifically, and water quality in general, requires development of a comprehensive database for lake water quality that is integrated with measures of lake ecological context and sufficient metadata for robust analysis.

We created a database called LAGOS-NE, the “lake multi-scaled geospatial and temporal database” for thousands of inland lakes in 17 of the most lake-rich states in the upper Midwest



**Figure 1:** Map of the study extent of LAGOS-NE. Map includes 17 states in the upper Midwest and Northeastern United States outlined in white and 51 101 lakes  $\geq 4$  ha shown as blue polygons. Some lakes extend beyond state borders and are included in the database if it was possible to delineate their watersheds. Watershed boundaries rather than state boundaries were used for all analyses of lakes, streams, and wetlands. The map is modified from Soranno et al. [17].

and Northeastern United States (Fig. 1). We avoided the problem of lack of metadata for the water quality data by contacting the original data providers for water quality data, asking for metadata, and only including data for which sufficient metadata were available. We addressed the problem of lack of ecological context data by creating our own database of lake ecological context. The detailed methods and approach for building this database have been published previously [17]; here we publish and describe the database for the 51 101 lakes and reservoirs  $\geq 4$  ha in the study area (1 800 000 km<sup>2</sup>).

We had 3 related motivations for developing this database: (i) to facilitate further development of our basic understanding of lake water quality at broad scales using water quality data on thousands of lakes collected over the last several decades (see [11, 17] for details); (ii) to build the capacity to apply this scientific understanding to environmental management and policy of inland waters; and (iii) to foster broad-scale research by designing an open-science database that is extensible for future uses and by making the data and methods publicly accessible.

LAGOS-NE comprises 3 data modules that, although integrated in the same database, were derived using different data sources and data integration methods, and thus must be version-controlled separately. LAGOS-NE<sub>LOCUS</sub> v1.01 includes lake location and physical characteristics based on an existing national-scale database of lake and streams in the United States for all lakes. LAGOS-NE<sub>GEO</sub> v1.05 includes measures of land, water, and air (ecological context) obtained from existing national-scale GIS (geographic information system) data sets and measured in multiple zones (delineated by different spatial classifications) around all lakes. This module also contains some

temporal data for climate, land use/cover, and atmospheric deposition variables. LAGOS-NE<sub>LIMNO</sub> v1.087.1 includes in situ measurements of lake water quality for a subset of the above lakes. These 87 data sets of lake water quality were obtained from a combination of sources including government, tribal agencies, university researchers, citizen scientists, and non-profit agencies. Samples were taken during any season of the year from the most recent decades, mostly from the late 1980s to 2012.

The largest challenge in building LAGOS-NE was the heterogeneity of the data set formats, variable conventions and units, and metadata, none of which were standardized. Many steps of data integration required manual input from experts in diverse fields and close collaboration among specialists in ecoinformatics, database design, freshwater ecology, and geography; all combined, the effort took 6 years and involved  $\sim 15$  individuals, spread across numerous institutions.

We designed the database using principles of open science so future users could ask new research questions by using the existing database or adding new data modules to the database. To ensure that users could do this, we documented the major steps of data set integration and carefully integrated metadata directly into the database itself, we emphasized data provenance, and we used a database versioning system. In this data paper, we make the following research products available: (i) data tables with the data that make up LAGOS-NE and an R package for accessing the data and integrating the tables; (ii) for each of the 87 water quality data sets, we provide the ecological metadata language (EML) metadata files that we authored after receiving the data, the data files that we processed to import into LAGOS-NE and the R-script that we wrote

**Table 1:** Summary statistics for LAGOS-NE study area

State	Area (km <sup>2</sup> )	Number of lakes (≥4 ha)	Mean annual temperature (°C)	Mean annual precipitation (mm)	% agricultural land	% urban land	% forested land	% wetland
Connecticut	12 878	763	9.7	1253	7.2	24.4	54.5	9.0
Illinois	145 920	2819	11.3	1005	68.9	11.9	15.0	1.7
Indiana	93 717	1874	11.2	1072	62.0	10.8	22.5	1.5
Iowa	145 736	903	9.1	881	78.0	7.5	6.9	1.9
Maine	84 123	2645	5.1	1149	3.7	3.5	66.9	12.1
Massachusetts	21 013	1698	8.9	1235	5.8	25.2	50.1	12.2
Michigan	150 489	6511	7.2	841	26.2	10.6	35.5	19.2
Minnesota	218 543	13 984	5.3	709	44.7	5.7	19.7	19.0
Missouri	180 537	1858	12.7	1100	50.7	7.0	36.6	2.1
New Hampshire	23 980	1109	6.5	1209	3.8	7.9	74.5	6.4
New Jersey	19 599	1143	11.8	1188	13.8	31.1	27.9	21.4
New York	126 070	4461	7.6	1094	21.9	9.3	54.1	7.2
Ohio	106 917	1279	10.6	1003	50.0	14.7	30.9	1.0
Pennsylvania	117 293	1755	9.3	1109	22.7	12.3	59.5	1.6
Rhode Island	2809	253	10.0	1246	4.9	29.5	44.6	13.6
Vermont	24 913	528	5.9	1176	13.3	5.5	70.0	4.7
Wisconsin	145 295	6009	6.6	831	36.7	7.5	35.5	13.7

This table includes the numbers of lakes and geophysical setting of each state and state averages for climate and the 4 major land use/cover types, which do not add up to 100% because we do not include all cover types. Temperature and precipitation data are 30-year climate norms (1981–2010; PRISM, <http://www.prism.oregonstate.edu/normals/>); land use/cover data are from the 2011 National Land Cover Database (NLCD; USGS, <http://www.mrlc.gov>). Note, border lakes are only counted in 1 state.

to process the data; and (iii) GIS coverages of the underlying freshwater geographic features (lakes, streams, and wetlands) that are linked to the data tables for GIS processing by researchers.

### Study Site: Midwest and Northeast US Lakes

We selected an area of the United States known to have large numbers of lakes, well-developed lake water quality sampling programs, and that spans diverse geographic conditions and thus gradients of ecological context (Table 1). Our study area of 17 US states includes 51 101 lakes ≥4 ha (Fig. 1). These states are in the north temperate climatic zone, which experiences cold winters and warm, humid summers. The study area includes part of the Interior Plains, Laurentian Uplands, Appalachian Highlands, and Atlantic Plain geological provinces, and thus encapsulates a range of geological ages, glacial histories, and topography. Land use/cover is highly variable, ranging from regions of intense agriculture in the corn belt that span portions of Minnesota, Wisconsin, Iowa, Missouri, Indiana, and Ohio, to predominantly forested or urban regions of the northeastern United States, including the states of Maine, New Hampshire, New Jersey, and parts of New York, and primarily forested regions of northern Minnesota, Wisconsin, and Michigan.

Although the majority of the data that we provide are for lakes ≥4 ha (see below for reasons for using this threshold), we do include some data on lakes ≥1 ha and <4 ha if data were available. Although there may be water quality data for some lakes in this smaller size range, ecological context variables are not available for these lakes.

### Overview of LAGOS-NE

LAGOS-NE includes some data on all lakes in a study area (above the minimum lake area threshold, which was 4 ha), which we call the “census” population of lakes. The census population of lakes is a critical feature of LAGOS-NE because it allows us to

characterize the ecological context of every lake in our study population and to identify whether the lakes for which we have water quality data are biased in any way. LAGOS-NE includes 3 main categories of variables: (i) variables that describe the physical characteristics and location of lakes themselves; (ii) variables that describe in situ water quality; and (iii) variables that describe a lake’s ecological context at multiple scales and across multiple dimensions (such as hydrology, geology, land use, climate, etc.) based on the principles of landscape limnology [12, 18–20]. Three factors dictated which data were included: past research and theory about the spatial and temporal controls of lake water quality, data availability and quality, and the time and resources necessary to compile, integrate, and process the original data. In other words, data that were especially time- and resource-intensive to collate, integrate, or process were given lowest priority and, in some cases, were not ultimately incorporated into the database.

There was a number of constraints for each of the categories of data that had to be considered. For creating the census population of lakes (i.e., their geospatial location, perimeter, and surface area), we relied on a single source of data (the 1:24 000 National Hydrography Dataset [NHD]) [21]. For the in situ water quality data, we incorporated data only if they were in a digitally accessible format such as a text or spreadsheet file. Finally, for the ecological context variables, we included only data for which we could obtain a GIS or raster coverage at the national or state scale for all 17 states.

We organized these 3 categories of data into database “modules” that had similar data types and sources so that we could develop procedures and set standards for each module (Fig. 2). The module structure also facilitates data reuse and extension by accommodating future data modules related to any other lake or ecological context feature.

The design of LAGOS-NE and the workflow for its construction have been described previously in detail [17]. In particular, the database design is based on the Consortium of Universities for the Advancement of Hydrologic Science, Inc. (CUAHSI),

## LAGOS-NE Modules

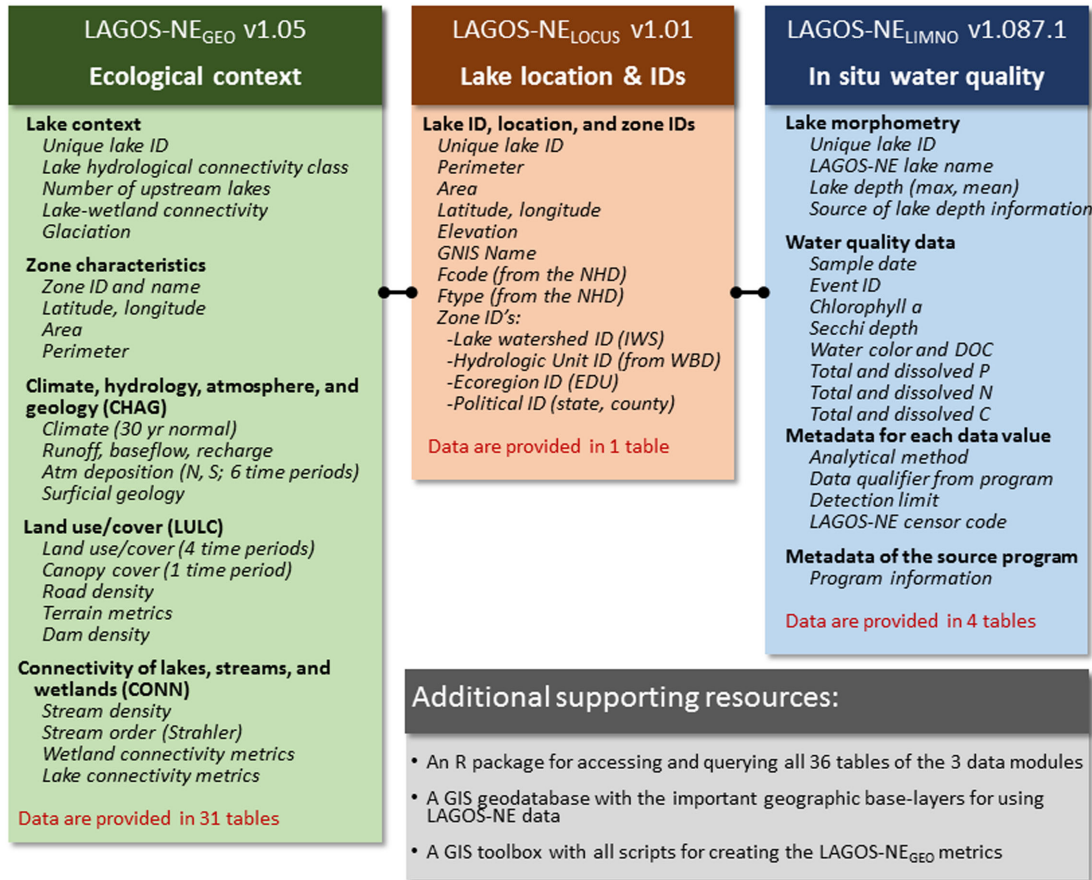


Figure 2: LAGOS-NE data modules and version numbers. The data modules and versions that are included in LAGOS-NE and are available with this paper include LAGOS-NE<sub>GEO</sub> v1.05, LAGOS-NE<sub>LOCUS</sub> v1.01 (note that in Soranno et al. [17], this module was called LAGOS-lakes), and LAGOS-NE<sub>LIMNO</sub> v1.087.1. We include descriptions of the types of data that are included in each module, with the major categories of variables the same as those describing the data tables in Additional file 1. The black connectors among the modules show that the modules are connected to each other through common unique identifiers through the LAGOS-NE<sub>LOCUS</sub> module (through the unique lake ID). P is phosphorus, N is nitrogen, C is carbon, S is sulfur, and atm is atmospheric. This figure is modified from Fig. 1 in Soranno et al. [17].

Community Observations Data Model (ODM; CUAHSI ODM) as described in Soranno et al. [17]. Here, we provide a brief overview. One important guiding principle in creating LAGOS-NE was to ensure data provenance, i.e., that we could trace the original source data through to the final LAGOS-NE database. Because each data module had different types of source data, we developed different procedures for data provenance for each module, described in Soranno et al. [17] and in this paper. The database model is based on ODM because it is a flexible data model (i.e., allows the incorporation of a wide range of types of data) that allows for the incorporation of controlled vocabulary and, importantly, allows for extensive documentation through a relational database structure of linked tables containing metadata [17]. The database was created and is maintained in PostgreSQL v9.1. However, for researchers to use the database for analysis and modeling, it is necessary to export the data into tables that can be processed by statistical packages or computer code. Therefore, we exported the data into a series of tables (of similar data) that are needed to conduct research on either the census population of lakes, the lakes for which there are water quality data, or some combination. These are the data files that have been used to conduct research on LAGOS-NE to date and that we make available in this paper (see Additional file 1 for a list of the tables and associated data that

we are making available). Further, we also make our GIS data sets available to facilitate geospatial analyses of lakes, streams, and wetlands used to create some of the major components of LAGOS-NE.

### Description of the LAGOS-NE<sub>LOCUS</sub> v1.01 data module

The LAGOS-NE<sub>LOCUS</sub> module includes data on the physical location, some features, and unique identifiers for all lakes in the study area  $\geq 1$  ha, which means this data file has information on 141 378 lakes. Note that, because we detected errors in the digitization of lakes between 1 and 4 ha, we have chosen to define our census population of lakes as only those  $\geq 4$  ha, but we still make data available for lakes smaller than 4 ha when available in this and the LAGOS-NE<sub>LIMNO</sub> data module. However, we recommend caution in analyses, interpretation, and inference for lakes  $< 4$  ha in this database that depend on NHD's spatial representation and detection of water bodies. The data in this module include lake unique identifiers, perimeter, area, latitude and longitude (which are typically the centroid of the lake or a central point that is within the lake boundary), GNIS name, and the zone IDs that the lake is located within (e.g., state, county, or hydrologic units). The GIS data sets that we also make available provide the

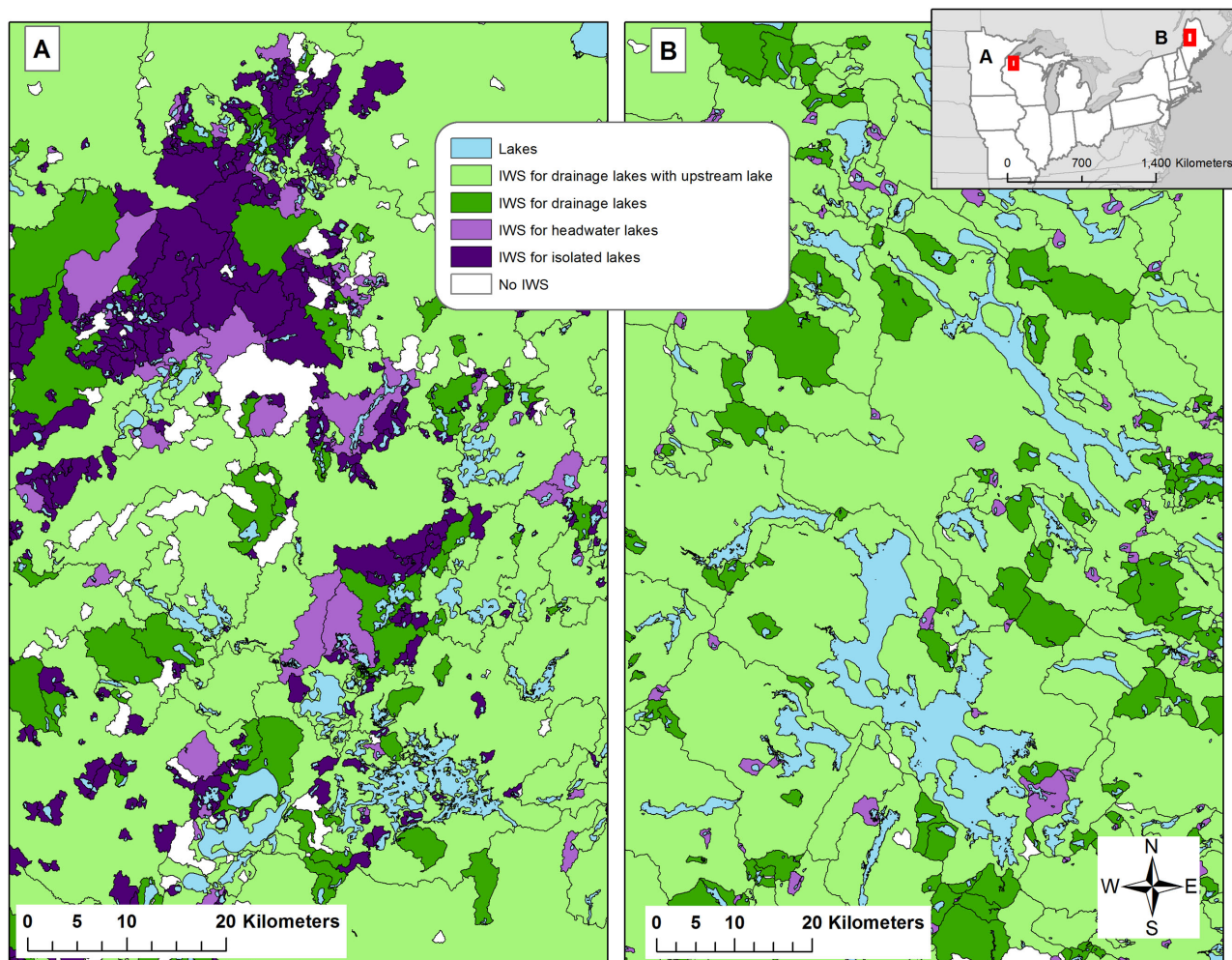


Figure 3: Examples of lake watersheds in LAGOS-NE. The watersheds are coded by the hydrologic class to which each lake belongs. Data are from the LAGOS-NE<sub>GEO</sub> v1.01 data module and the GIS data coverages.

lake polygon features associated with this module, as well as coverages for lake watersheds, streams, wetlands, spatial classifications, and glaciation history.

### Definition of lakes

We defined lakes previously in Soranno et al. [17] as follows. A “lake” in LAGOS-NE is a perennial body of relatively still water. We include lakes and reservoirs that range from being completely natural to highly modified: lake basins can be entirely natural, modified natural (i.e., a water control structure on a natural lake), or a fully impounded stream or river (i.e., a reservoir). We explicitly exclude sewage treatment ponds, aquaculture ponds, and detention ponds that are known to contain basins that are entirely artificial and were built for high-intensity human use. In addition, due to their unusual nature and size, we do not include the 5 Laurentian Great Lakes in our database. This definition of “lake” for LAGOS-NE has been developed only for the purpose of this database and its applications (e.g., to answer questions about lake water quality). The intent of LAGOS-NE is not to document and measure the total number of water bodies in our study area, although we are able to perform this calculation for lakes  $\geq 4$  ha with an acceptable level of uncertainty (see below).

### Definition of lake watersheds

We calculated lake watersheds as “inter-lake watersheds” (IWS), defined as the area of land draining directly into the lake as well as the area that drains into upstream-connected streams and lakes  $< 10$  ha (Fig. 3). We defined lake watersheds this way to define the drainage basin of lakes that includes connected streams and their drainage basins. However, because research has shown that large upstream lakes can trap nutrients flowing into them, these large lakes can block the transport of nutrients that originate upstream from them to downstream lakes in a connected lake chain (e.g., [22]). Therefore, to calculate a drainage basin for a lake with large upstream connected lakes, we did not include the drainage basins of upstream lakes  $> 10$  ha. See Soranno et al. [17] for full details on how lake IWSs were calculated and the section on LAGOS-NE<sub>GEO</sub> for further details.

### Lakes near and beyond the state borders

For some of our analyses, we delineated boundaries in other ways than political boundaries that were more ecologically relevant, which resulted in the inclusion of some lakes outside of the exact 17-state border. This fact allowed us to include more in situ data collected by state and citizen sampling programs that

do not always follow strict state borders and may include lakes that are outside of state lines. Although most of these border lakes have hydrologic (i.e., lake connectivity measures) and topographic (i.e., lake watershed delineations) calculations or water quality data, some measures of ecological context may be missing. For example, for lakes in Canada, we were not able to estimate any data that relied on national data sets that stopped at the Canadian border; one exception is the NHD, which extends into Canada to retain hydrologic boundaries.

### Data sources of the LAGOS-NE<sub>LOCUS</sub> module

Detailed information on data sources are found in Additional file 5 in Soranno et al. [17]. Briefly, the data source for lakes and streams in the 17-state area was the NHD [21]. The hydrologic boundaries (i.e., for 3 of the spatial classifications, HUC12, HUC8, HUC4) came from the Watershed Boundary Dataset (WBD) [23]. In addition, we used the digital raster data set of elevation for watershed delineation from the National Elevation Dataset [24]. All download dates for these data sources are provided in Additional file 5 in the above citation.

### Data-integration methods of the LAGOS-NE<sub>LOCUS</sub> module

All methods to create this module are described in Soranno et al. [17]. The most challenging and time-consuming part of building this module was connecting the sampling locations from the lake water quality data sets (each contained different types of unique identifiers, and sometimes only lake names) to a georeferenced location in the NHD. When data providers included the lake latitude and longitude, we were able to mostly automate the procedure. Nevertheless, even when coordinates were available, there were many cases where the latitude and longitude did not intersect with the NHD lake polygon boundary, requiring manual interpretation.

### Quality Control of the LAGOS-NE<sub>LOCUS</sub> module

The full description of error analysis for this module is described in Soranno et al. [17]. However, here we briefly describe our efforts to determine the minimum area of a lake that we could confidently represent using the NHD (further details located in Additional file 9 in Soranno et al. [17]). Although the NHD is a national data set, it is updated and edited regionally (often at the state level) by local practitioners familiar with each study region. As a result, there are regional differences in the resolution and digitization of water bodies, particularly for small water bodies, making it difficult to quantify or document even nominal error rates, or rather, the minimum lake size that is well-represented in the NHD. It has been documented previously that the NHD may not successfully identify small water bodies due to a variety of reasons including the resolution of the original underlying data of the NHD database, errors in digitization, and hydrologic changes since the time of map creation (e.g., [25, 26]). Because of these documented issues, some programs have set minimum lake area cutoffs for sampling lakes. Most notable is the EPA-National Lakes Assessment of 2007, which chose a minimum size of 4 ha, although a smaller size cutoff was chosen for the EPA-National Lakes Assessment of 2012 [27]. To determine an appropriate size cutoff for our purposes, we conducted an analysis to identify the lakes that are best represented by the NHD across the LAGOS-NE study area.

We selected 4 states (WI, MI, IA, ME) in which to evaluate error rates of water body identification for lakes  $\geq 1$  ha and 7 states (WI, MI, IA, ME, MO, NH, OH) in which to evaluate error rates for lakes  $\geq 4$  ha. We randomly selected three 100-km<sup>2</sup> rectangles from each state then compared the number of lakes occurring in the NHD GIS coverage with the number of lakes in the best available aerial imagery from a range of sources to calculate the percentage of lakes missing from the NHD. The average percentage of lakes missing from the NHD was 58% for the  $\geq 1$  ha 4-state test and 13% for the  $\geq 4$  ha 7-state test. Because an average of 87% of lakes  $\geq 4$  ha that are present in high-resolution aerial imagery are also present in the NHD, we chose this surface area as our cut-off and accepted this error rate.

### Data in the LAGOS-NE<sub>LOCUS</sub> module

Figure 1 shows the census population of all lakes  $\geq 4$  ha in the 17-state area, including border areas beyond the 17-state boundary. As expected, the lakes are not evenly distributed, with higher densities in the northern parts of the study area. For those lakes with known lake depth (9808 lakes with maximum depth values, and 4090 lakes with mean depth values), there is little regional pattern of lake depth; shallow and deep lakes are found throughout the study area (see [28] for further details). Watershed size varies greatly across the study extent, reflecting the wide range of different lake hydrologic types and connections to upstream water bodies (Fig. 3). In fact, the proportion of lakes in different lake hydrologic connectivity classes varies regionally across our study extent (Table 2) (see [29] for further details).

### Description of the LAGOS-NE<sub>LIMNO</sub> v1.087.1 Data Module

The LAGOS-NE<sub>LIMNO</sub> module includes in situ measurements of lake water quality. We included variables that are most commonly measured by state agencies and researchers for studying eutrophication (water quality data and metadata, including chlorophyll a, Secchi depth, water color, DOC, total and dissolved phosphorus [P], nitrogen [N], and carbon [C]) (Fig. 2). For each water quality data value, we also include metadata as additional columns in the exported data table (metadata including analytical method, data qualifier from the program, detection limit [when available], and the LAGOS-NE sensor code) (Fig. 2), including the analytical methods, qualifiers with data flags from the original program (*qual*, which is not standardized for LAGOS-NE), detection limits (if available), and standardized sensor codes from our quality control procedures (*sensorcode*, standardized for LAGOS-NE). Finally, we include documentation about each source program that is linked to each data value.

### Data sources of the LAGOS-NE<sub>LIMNO</sub> module

We acquired individual water quality data sets for LAGOS-NE<sub>LIMNO</sub> by contacting individuals at each of the 17 state and 5 tribal agencies. These contacts helped us to identify the state agency-collected data set required by the Clean Water Act that was most likely to be in the public domain. In this way, we were able to acquire at least 1 (and typically more) data set from each of the 17 states. Because state and tribal agencies vary in sampling approach and intensity (see below for details), we sought to supplement these data sets with other known sources of water quality data, including university researchers, federal agencies, and non-profit groups, to integrate into the LAGOS-NE<sub>LIMNO</sub> module. The full list of data sources acquired is in Soranno



**Table 2:** Numbers of lakes in each state by lake hydrologic class

State	Lakes $\geq 4$ ha (#)	Isolated Lakes (#)	Headwater lakes (#)	Drainage lakes (#)	Drainage lakes with upstream lakes (#)
Connecticut	770	40	119	424	187
Illinois	2831	1417	279	952	183
Indiana	1883	760	244	697	182
Iowa	915	339	87	402	87
Maine	2661	94	619	1211	737
Massachusetts	1716	210	269	751	486
Michigan	6531	2649	1087	1672	1123
Minnesota	14 031	6609	1894	2673	2855
Missouri	1865	435	179	1113	138
New Hampshire	1118	70	224	581	243
New Jersey	1148	219	129	521	279
New York	4477	629	1210	1915	723
Ohio	1282	543	105	520	114
Pennsylvania	1757	316	397	840	204
Rhode Island	266	35	40	115	76
Vermont	531	14	74	364	79
Wisconsin	6026	2982	823	1236	985
Total	49 808	17 361	7779	15 987	8681

The number of lakes  $\geq 4$  ha in each of the lake hydrologic classes by state, as well as the total numbers of lakes by hydrologic class calculated for the study extent. Note, in this table, lakes are counted for each state in which they occur (i.e., lakes that straddle 2 states are counted in both states).

et al. [17] in Additional file 17; however, we incorporated a subset of these data sets in LAGOS-NE<sub>LIMNO</sub> v1.087.1 (the data file LAGOSNE'source.program\_10871.csv contains the list of sources for this version of LAGOS-NE).

### Data integration methods of the LAGOS-NE<sub>LIMNO</sub> module

All methods to create this module are described in Soranno et al. [17]. Briefly, for each data set acquired, we authored LAGOS-NE metadata in EML to aid in data provenance (included in this paper). We also incorporated key metadata features (e.g., methods used, sensor codes, if applicable), and sampling program information) into the database so that future users could easily identify these important attributes. Because each data set was unique in structure, file format, and naming conventions, we manually processed each data set and its metadata so that they could be translated into the standard LAGOS-NE vocabulary and data model. Although labor-intensive, we created customized R scripts to process and load each data set separately (included in this data paper).

### Quality control of the LAGOS-NE<sub>LIMNO</sub> module

The full description of our quality assurance/quality control (QAQC) procedures for this module is described in Additional file 2. Here, we provide a brief overview of our approach. Our goal for this effort was to identify egregiously high values and values that might be too low, both defined below. Note that our quality control procedures were not designed to identify statistical outliers, which individual users are expected to perform themselves because such analyses depend on the subsequent statistical analysis of each user. There were 3 major phases in the QAQC procedure for LAGOS-NE<sub>LIMNO</sub>. Phases I and II were designed to identify the egregious values that we defined as those that (i) did not make ecological sense, (ii) were far beyond what has been detected in previous studies, (iii) were not technically feasible (e.g., SRP > TP), or (iv) were a result of a data or file corruption or error in the data loading stage. For these

egregious values, we explored the issues that might be underlying the values and removed them from the LAGOS-NE<sub>LIMNO</sub> data export provided in this data paper because we had sufficient evidence that they were not scientifically valid data values. We were very conservative in these assessments to avoid removing data values that were high, yet still valid. Phase III was designed to identify and flag values that seemed to be lower than analytically possible (i.e., below detection limits) when there were sufficient metadata; however, note that these data are still provided in this data paper because it is not appropriate to remove data that are below detection when those data could be valid.

For all versions of LAGOS-NE<sub>LIMNO</sub>, phases I and II are conducted on the entire cumulative data set to leverage as large of a sample size as possible to detect problem values. In other words, because many of the QAQC analyses outlined here make use of all information from an individual lake or variable, incorporating new data may result in a better assessment of the data than when there are fewer data. Thus, for each new version of LAGOS-NE<sub>LIMNO</sub>, new decisions are made about egregious values. In this data paper, we describe the procedures for assessing all major versions of LAGOS-NE<sub>LIMNO</sub>, but we present the results only for this version of LAGOS-NE<sub>LIMNO</sub> (v1.087.1).

Because there are few accepted practices for conducting such quality control on a large, integrated database, we created our own procedures for phases I and II by creating tests to identify egregious values that leverage a large, integrated database with multiple measures of water quality and well-established expected relationships among variables. The database that we used to identify egregious values was based on data in the full LAGOS-NE<sub>LIMNO</sub> database for samples taken from all lake depths provided by the source data sets (note, our data exports in this data paper are only for epilimnetic or surface samples). While the quality control procedures that we implemented here were designed to help resolve the large and egregious errors in a combined data set such as this, there are likely additional extreme values in the database due to the size and heterogeneity of the data. Users may want to check for additional issues in the data values specific to their intended analyses.

## Data in the LAGOS-NE<sub>LIMNO</sub> module

All data in LAGOS-NE<sub>LIMNO</sub> v1.087.1 are from samples that we identified as being collected from either the lake surface or the epilimnion (the well-mixed surface layer of a thermally stratified lake during the period of stratification). Because we did not have lake temperature data to quantify the exact epilimnion depth in all lakes, we used information from the source data sets to either determine epilimnion depth or to select data from only the top water layers. Although we received data from different depths in lakes, the majority of the samples were from the surface or epilimnion. The database includes samples from any season of the year. However, most of the published analyses to date have focused on the summer stratified period.

Lakes are not sampled the same way by all individuals, groups, or agencies; there are differences in the variables measured, the frequency and timing of sampling, and the proportion of lakes sampled. For example, for total phosphorus, the 4 states with the largest number of unique lakes with at least 1 value for total phosphorus per state include Wisconsin (1920 lakes), Minnesota (1588), New York (1289), and Michigan (1109) (Table 3). However, the states with the highest proportion of their lakes with total phosphorus samples are the smaller states with fewer lakes, such as New Hampshire (64%), Vermont (58%), and Rhode Island (42%). Notably, there are some states with intermediate numbers of lakes that still have quite large percentages of their lakes with total phosphorus values, including Maine (35% of 2645 lakes), Wisconsin (32% of 6009 lakes), and New York (29% of the 4461 lakes).

The most commonly measured variable in LAGOS-NE<sub>LIMNO</sub> is water clarity, measured as Secchi depth (a relatively easy and cost-effective measure of water quality), with 897 724 measurements taken from 12 034 unique lakes in the 17 states from mostly the mid 1980s to 2011 (Table 3). The second and third most sampled measures of water quality are chlorophyll a and total phosphorus, respectively. Although it appears that total nitrogen is sampled far less frequently than total phosphorus, some labs measure total nitrogen directly and report that single value, whereas other labs measure the constituents that make up total nitrogen (total Kjeldahl nitrogen and nitrate+nitrite) and sum them together to calculate total nitrogen. All of our analyses conducted on total nitrogen have used such calculated and measured values of nitrogen together, which increase the sample sizes for total nitrogen markedly.

Most of our data came from state agencies, either alone or as part of joint programs with citizen scientists or university researchers (Table 4), which highlights the importance of citizen science programs for monitoring lake water quality in this lake-rich area of the United States.

Using the 3 most sampled variables in the data set (Secchi depth, chlorophyll concentration, and total phosphorus), we found that larger lakes were more likely to be sampled for water quality than smaller lakes (Fig. 4). This result was expected given the economic and recreational interest in larger lakes, including easier public access. Previous research has already documented this basic pattern in 6 of the states included in LAGOS-NE [30]. Across all states, almost 80% of lakes >400 ha have water quality data.

Lakes are also unevenly sampled through time, depending on the variable (Fig. 5). Some programs' focus is on long-term monitoring, whereas others are short-term initiatives. Typically, long-term monitoring programs are localized to a few lakes, although there are exceptions (e.g., monitoring for acid rain in the northeastern United States in the 1980s-present has resulted in

good temporal and spatial coverage for some variables through time and space) [31].

## Description of the LAGOS-NE<sub>GEO</sub> v1.05 Data Module

The LAGOS-NE<sub>GEO</sub> module includes information on the ecological context of the census lakes, their watersheds, and their regions. The information provided in the data tables for this module is organized into 3 main themes in which data are exported into individual tables: CHAG—climate, hydrology, atmospheric deposition of nitrogen and sulfur, and surficial geology; LULC—land use/cover, canopy cover, terrain metrics, and dam density; and CONN—lake, stream, and wetland abundance and connectivity measures (Fig. 2). We also provide the GIS coverages that include some of the underlying data for this module, including lake polygons and their hydrologic classifications, defined in Soranno et al. [17]; wetland polygons and their classification; streams as a line coverage and their classification by stream order; the zones used for this study (state and county, hydrologic units [at the 4, 8, and 12 scales]) [32]; and lake watersheds (IWS). We also include boundaries of US states and Canadian provinces for mapping.

### Data sources of the LAGOS-NE<sub>GEO</sub> module

Detailed information on data sources are found in Additional file 5 in Soranno et al. [17]. Almost all data sources for this module are from national-scale data sets and thus use standardized methods throughout the study extent.

### Data integration methods of the LAGOS-NE<sub>GEO</sub> module

All methods to create this module are described in Additional files 5, 7, 8, 13, and 14 in Soranno et al. [17]. Briefly, we calculated the metrics for this module that describe the ecological context surrounding lakes by developing project-specific GIS tools in the ArcGIS environment, which are referred to as the LAGOS GIS Toolbox [33]. The toolbox outputs multiple individual data tables of calculated values organized by the above 3 data themes that are then imported into LAGOS-NE<sub>GEO</sub> for different spatial classifications, including values calculated at the level of the individual lake, 100-meter and 500-meter buffers around each lake, the lake IWS, states and counties, hydrologic units, and ecological drainage units (an ecoregion spatial classification). The unique identifiers for this data module are the zone IDs for each spatial classification for which we calculate these metrics. In other words, we calculate land use around a lake in each of the zones of the many spatial classifications in LAGOS-NE. However, the data are exported into individual tables by spatial classification. Therefore, there are different numbers of rows in each table; for example, there are 51 101 rows for the land use metrics calculated for the 100-meter lake buffer because there are 51 101 lakes that have a 100-meter buffer area, but only 17 rows for the land use metrics calculated for the state spatial classification.

### Quality control of the LAGOS-NE<sub>GEO</sub> module

The full description of error analysis for this module is described in Additional file 14 in Soranno et al. [17]. The quality control procedures for this module included procedures to identify possible errors or improbable values as a result of the extensive automated GIS data processing that creates the LAGOS-NE<sub>GEO</sub> data tables and to correct those problems. We assumed that the

Table 3: Summary of the water quality variables and the number of values per variable by state

State	Number of lakes ( $\geq 4$ ha)	Variable	Total phosporus	Secchi depth	Chlorophyll a	True color	Apparent color	Dissolved organic carbon	Total nitrogen	Total Kjeldahl nitrogen	Nitrate + nitrite
Connecticut	763	# of samples	1294	1943	1160	53	0	74	853	55	397
		# of sampled Lakes	143	168	149	37	0	49	99	26	81
Illinois	2819	Sample years	1972–2010	1937–2010	1937–2013	1984–2007	n/a	1984–2007	1973–2010	1999–2009	1976–2010
		# of sampled lakes	2816	2317	1438	20	0	20	43	1526	2351
Indiana	1874	# of sampled lakes	191	185	167	17	0	17	18	155	188
		Sample years	1999–2011	1999–2011	2000–2011	2007	n/a	2007	2001–2009	1999–2006	1999–2009
Iowa	903	# of samples	1232	1303	909	57	0	57	57	1183	1237
		# of sampled lakes	341	340	320	51	0	51	51	322	341
Maine	2645	Sample years	1988–2010	1986–2010	1990–2009	2007	n/a	2007	2007	1988–2009	1988–2009
		# of sampled lakes	2873	2836	2711	18	0	18	2244	6	2229
Massachusetts	1698	# of samples	111	111	103	12	0	16	111	1	111
		Sample years	1997–2011	1997–2011	1997–2011	2007	n/a	2007	2001–2011	2008–2009	2001–2011
Michigan	6511	# of samples	17314	83472	12480	1927	1676	3321	1260	8	1577
		# of sampled lakes	933	1047	793	601	466	848	461	3	347
Minnesota	13984	Sample years	1971–2011	1952–2011	1974–2011	1983–2011	1972–2011	1984–2011	1995–2011	1978–1993	1978–2011
		# of samples	570	760	326	277	228	300	300	69	69
Missouri	1858	# of sampled lakes	211	249	122	122	89	140	37	4	132
		Sample years	1978–2013	1978–2010	1986–2010	1984–2013	1978–2010	1984–2010	2000–2010	1978–2013	1978–2013
New Hampshire	1109	# of samples	10143	95283	12243	1811	69	987	749	2651	4850
		# of sampled lakes	1109	1233	862	836	69	353	200	713	948
New Jersey	1143	Sample years	1965–2013	1925–2013	1959–2013	1973–2010	2002–2003	1984–2013	1959–2011	1980–2010	1973–2012
		# of samples	1588	497646	81925	406	6683	3382	7717	619	43054
New York	1143	# of sampled lakes	1588	4118	2755	253	1368	811	619	2018	1522
		Sample years	1944–2011	1938–2012	1970–2012	1981–2009	1949–2011	1984–2012	1945–2012	1944–2012	1945–2012
Ohio	1858	# of samples	11619	11794	11578	27	0	27	11340	0	27
		# of sampled lakes	208	207	201	23	0	23	207	0	23
Pennsylvania	1109	Sample years	1978–2013	1978–2013	1978–2013	2007	n/a	2007	1978–2013	n/a	2007
		# of samples	9289	2958	154	237	3044	390	390	22	1209
Virginia	1143	# of sampled lakes	710	618	21	111	603	143	17	535	704
		Sample years	1975–2013	1975–2011	1983–2012	1984–2010	1975–2010	1984–2010	2004–2010	2004–2010	1975–1994
Washington	1143	# of samples	421	461	446	27	0	44	10	443	472
		# of sampled lakes	175	174	157	25	0	36	8	157	175
Wisconsin	1858	Sample years	1984–2009	1984–2009	2005–2009	1984–2007	n/a	1984–2007	2007	2005–2009	1984–2009
		# of samples	175	174	157	25	0	36	8	157	175

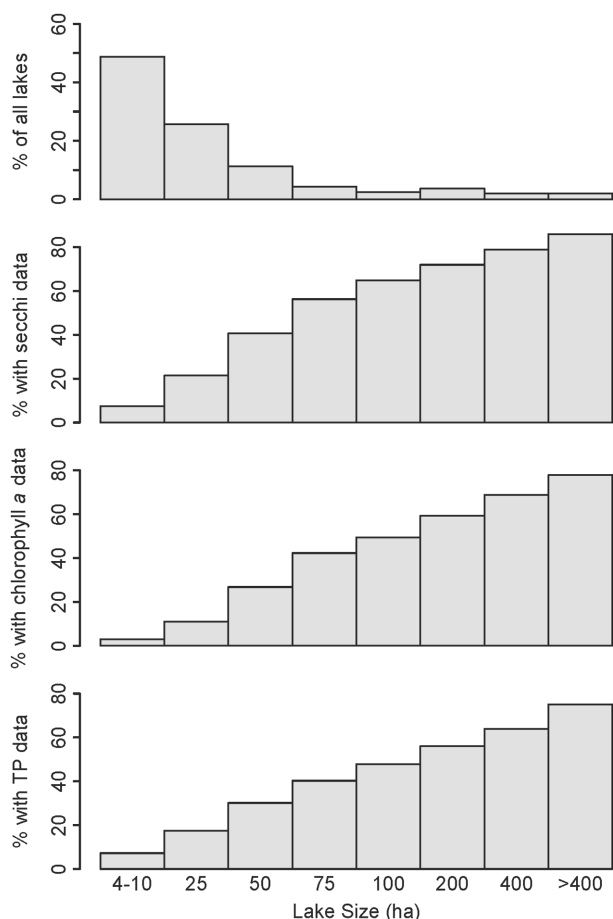
Table 3: Continued

State	Number of lakes ( $\geq 4$ ha)	Variable	Total phosphorus	Secchi depth	Chlorophyll a	True color	Apparent color	Dissolved organic carbon	Total nitrogen	Total Kjeldahl nitrogen	Nitrate + nitrite
New York	4461	# of samples	21356	21235	21000	27297	2287	13036	8259	944	27796
		# of sampled lakes	1289	693	545	1421	47	1158	258	279	1279
Ohio	1279	Sample years	1975-2012	1975-2012	1975-2012	1981-2012	1984-2011	1982-2011	1990-2012	1981-2010	1975-2012
		# of samples	377	1868	1912	20	0	220	1873	0	447
Pennsylvania	1755	# of sampled lakes	144	144	137	19	0	44	145	0	40
		Sample years	2006-2007	1992-2010	1992-2010	2007	n/a	2006-2010	1994-2010	n/a	1993-2007
Rhode Island	253	# of samples	1170	924	971	163	0	160	638	16	290
		# of sampled lakes	263	260	160	124	0	124	167	2	147
Vermont	528	Sample years	1980-2011	1984-2011	1980-2011	1984-2008	n/a	1984-2007	1997-2011	1985-2010	1980-2010
		# of samples	3325	18211	12195	51	6	65	2582	0	2100
Wisconsin	6009	# of sampled lakes	106	107	102	27	1	32	99	0	102
		Sample years	1984-2010	1984-2010	1986-2010	1984-2007	2003-2010	1984-2010	1992-2010	n/a	1984-2010
TOTAL	49592	# of samples	13906	23894	15273	1774	1542	982	8	194	2271
		# of sampled lakes	307	301	249	94	82	83	8	2	116
		Sample years	1977-2010	1977-2010	1977-2010	1981-2010	1979-2010	1984-2010	2007	1979-1994	1977-2010
		# of samples	45973	130819	26068	4599	174	4029	1932	9596	9417
		# of sampled lakes	1920	2079	1024	1281	1	671	180	1160	1216
		Sample years	1933-2013	1948-2013	1933-2013	1974-2013	1976-1998	1977-2013	1986-2010	1933-2013	1965-2013
		# of samples	154652	897724	202789	38764	15709	27112	39656	60954	65982
		# of sampled lakes	9749	12034	7867	5054	2726	4599	2685	5377	7472

We include the number of individual values (representing an individual sampling event), the number of unique lakes for which there is at least 1 data value, and the earliest and most recent year of sampling, all recorded by state and variable from any time period. Additional variables in LAGOS-NE<sub>LIMNO</sub> v1.087.1 not included in this table, which have relatively low sample sizes, include dissolved Kjeldahl nitrogen, ammonium, nitrite, soluble reactive phosphorus, total dissolved phosphorus, total organic carbon, and total organic nitrogen.

Table 4: The number of data sets, data values, and lakes from the different types of sampling programs in LAGOS-NE v1.087.1

Program Type	Number of data sets	Number of lakes ( $\geq 4$ ha)	Total phosphorus	Secchi depth	Chl. a	True color	Apparent color	Dissolved organic carbon	Total nitrogen	Total Kjeldahl nitrogen	Nitrate + nitrite
Federal agency	3	17	419	527	324	229	173	215	335	6	30
		# of values									
		# of unique lakes	17	17	17	13	15	14	16	1	9
Federal agency/university	2	2	-	799	-	-	-	-	-	-	-
		# of values									
		# of unique lakes	-	2	-	-	-	-	-	-	-
LTER	3	9	2346	3529	2567	-	-	1872	1612	507	2396
		# of values									
		# of unique lakes	9	9	5	-	-	9	9	4	9
National survey program	5	2244	2320	2595	243	3689	703	4714	431	-	4204
		# of values									
		# of unique lakes	1863	1891	171	13	142	2235	398	-	1997
Non-profit agency	4	44	1326	4798	2678	-	-	-	214	9	908
		# of values									
		# of unique lakes	44	41	28	-	-	-	39	1	44
State agency	33	4264	34 348	42 888	29 993	16 240	5010	14 528	5359	7220	25 684
		# of values									
		# of unique lakes	3914	3186	2309	2092	776	1191	634	1991	3216
State agency/citizen monitoring	11	7039	79 390	645 650	124 766	18 010	8630	3195	18 610	52 995	27 826
		# of values									
		# of unique lakes	3955	6629	4341	1111	1508	786	772	3476	2782
State agency/univ/citizen monitoring	4	1835	31 809	194 177	37 993	439	1171	1519	10 844	-	2112
		# of values									
		# of unique lakes	1439	1812	1253	302	393	574	712	-	99
Tribal agency	5	46	911	145	905	3	-	357	411	277	463
		# of values									
		# of unique lakes	33	3	32	3	-	11	18	5	17
University	17	535	2273	4412	3939	172	69	723	2275	-	2397
		# of values									
		# of unique lakes	326	500	415	151	69	318	396	-	171



**Figure 4:** Percentage of lakes by lake area with water quality data. Percentage of census lakes in each lake size bin (top panel) compared with the percentage of census lakes for which there are limnological data for Secchi (second panel), chlorophyll a (third panel), and total phosphorus (TP; bottom panel).

original data layers had already gone through extensive quality control by the originators of the data sets. We defined errors and improbable values to be: (i) values that did not make ecological sense; (ii) values that were well beyond what has been observed in previous studies; (iii) values that are not technically feasible; or (iv) null values that indicate an absence of data, when in fact data exist based on the input data coverages. Note, it was not our intention to remove statistical outliers that may or may not be real/true values. Rather, we conducted procedures on each exported table that included verifying column headers and units, mapping the exported data to evaluate mapping extent and boundary issues using visual inspection, mapping the data distributions of each value, identifying values that were missing or zero, plotting distributions of the data, ensuring that proportions summed to 100 where relevant, and inspecting univariate plots of metrics that are known to be related (e.g., % urban land use vs % impervious surface).

### Data in the LAGOS-NE<sub>GEO</sub> module

This module contains the largest amount of data of any of the modules. For example, Fig. 6 shows the wide range of ecological context for the LAGOS-NE study area calculated for 3 different spatial classifications. For those variables that are measured coarsely (e.g., baseflow, runoff, atmospheric deposition,

geology), we calculated variables for only the broader spatial classifications. For example, we did not calculate baseflow for spatial classifications finer than HUC12 because the underlying data for baseflow are estimated on a zone generally coarser than the area of a lake watershed.

### Research to Date Using LAGOS-NE

Prior versions of this database have supported numerous peer-reviewed publications to date. In particular, LAGOS-NE is ideally suited for studying the local to regional controls of water quality through both space and time because of the large number of lakes with in situ water-quality measurements and their wide gradients of ecological context. The lake census data set also makes it possible to quantify the types of biases present in the data set to assess the potential influence of uneven sampling efforts on results across both space and time. Below, we describe the types of research questions that have been and are being addressed using LAGOS-NE, organized according to 3 main topics related to studying water quality across space and time in thousands of lakes. We have published 10 articles using portions of this database, and 13 articles are in review or preparation presently.

### Methods and database development for macrosystems ecology

Several of our lines of research have required the development of novel methods and the application of existing methods in novel ways. Much of the impetus for this work on methods and database development has been driven by 2 needs. The first was to further develop the database—i.e., creating derived and predicted data as a new data product that is publicly accessible (e.g., [28]). The second was to better understand the spatial and temporal distribution of data contained in LAGOS-NE and to further our understanding of important ecological attributes of lakes across multiple spatial scales. These 2 needs are not mutually exclusive—analyses that have helped contribute data to LAGOS-NE have also addressed important ecological questions.

Three data gaps were identified early during database development, including (i) a lack of lake depth information (lake depth drives many in-lake processes), (ii) the need to develop a flexible method for creating ecological regions from multi-themed mapped data, which are often used in macroscale research to account for broad-scale patterns and processes, and (iii) the need for developing ways to measure freshwater connectivity to account for the transport and processing of materials in lakes at broad scales. For the first gap, Oliver et al. [28] used a linear mixed model to predict lake depth for lakes where in situ measurements were lacking, allowing the relationship between surface area and lake depth to vary by region because of the strong regional differences in this relationship. Predictions in some regions were far better than other regions, potentially due to differences in underlying geomorphology. To address the second gap, Yuan et al. [34] developed a novel spatially constrained spectral clustering algorithm that balances geospatial homogeneity and region contiguity to delineate ecological regions. Cheruvilil et al. [35] have since applied this clustering algorithm across the 17-state study region and tested the ability of newly developed regions to capture variation in lake nutrients and water clarity. Finally, to address the third gap, Fergus et al. [29] developed approaches for determining freshwater connectivity of lakes, streams, and wetlands across broad spatial

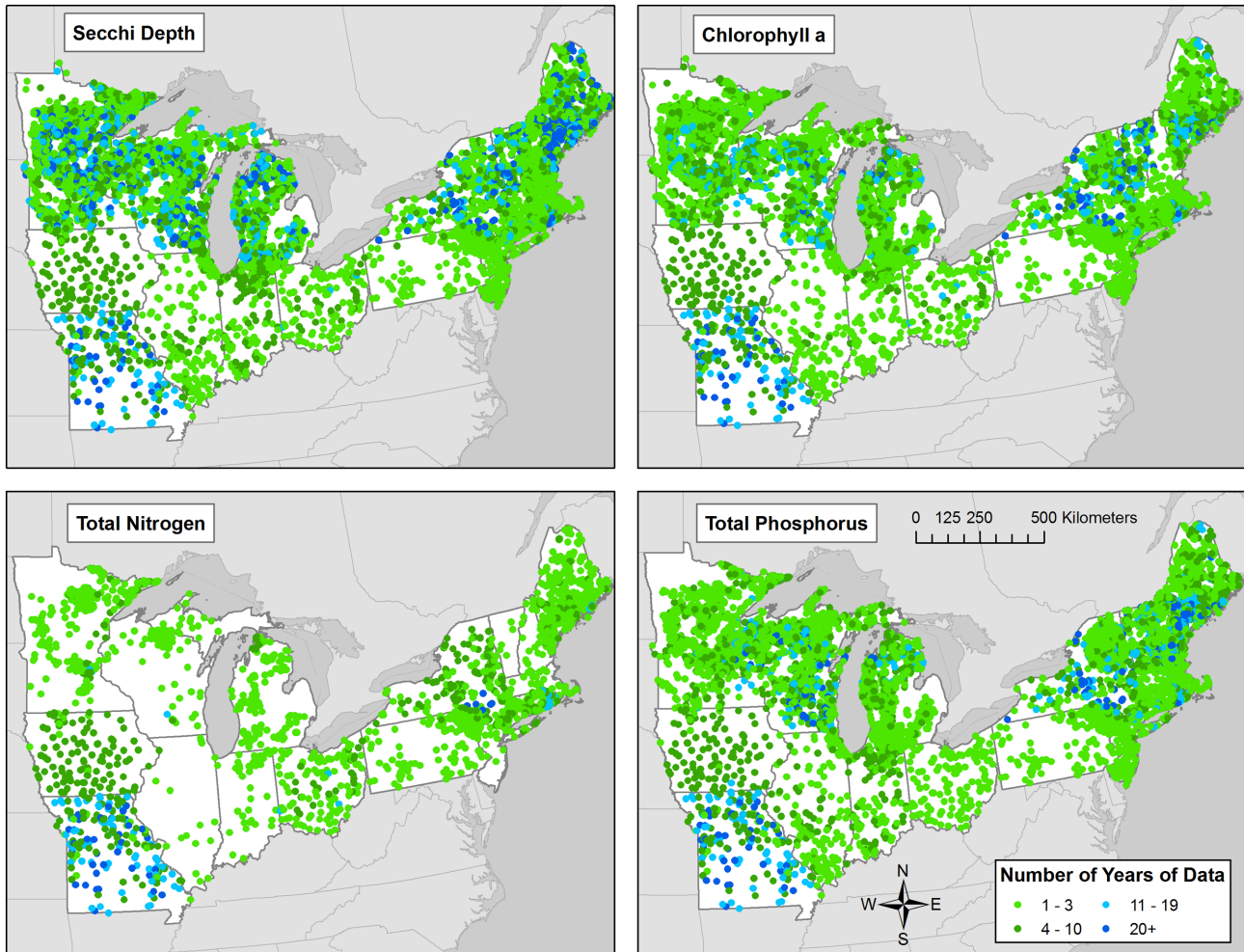


Figure 5: The number of years of water quality data by lake. The number of years for which at least 1 sample is taken during the summer stratified season (15 June to 15 September) for Secchi depth in meters, total phosphorus in  $\mu\text{g/L}$ , total nitrogen in  $\mu\text{g/L}$  (includes both measured and calculated values), and chlorophyll a in  $\mu\text{g/L}$ .

extents. The resulting freshwater metrics and analysis provide insight into the spatial distribution of surface water connectivity types across the LAGOS-NE study area and provide LAGOS-NE users with novel metrics of connectivity for use in future research.

A further challenge in large, integrated databases such as LAGOS-NE is the well-known problem with data derived from analytical methods related to the issue of detection limits [36]. Stow et al. (personal communication) studied the in situ concentrations that were too low to be quantified by standard analytical practices—measurements that are termed left-censored or below a detection limit of an analytical method. Unfortunately, detection limits were only sometimes reported (although we do include those data in LAGOS-NE<sub>LIMNO</sub> where available). In some cases, low values were flagged as being censored, with an explanation as to the reason for censoring the data value, but in other cases the reason for censoring was not clear. In some instances, patterns in the data suggested that ad hoc substitutions for censored observations may have occurred without clear documentation. Stow et al. (personal communication) describe a statistical approach that can be used to accommodate left-censored data during macroscale statistical analyses. This work also led to refining how censored observations were reported in LAGOS-NE, which has been incorporated into all later versions of LAGOS-NE<sub>LIMNO</sub>, including v1.087.1.

Lake water quality is affected by many ecological context features, such as lake physical characteristics, land cover, land use, and climate. The relationship between these features and the water quality measurements is not always linear. In addition, the data tend to be noisy and often contain missing values, which makes it challenging to fit effective statistical models. To overcome these challenges, Yuan et al. [37] developed a novel algorithm for learning non-linear features to predict lake water quality. The algorithm also enables the missing values to be imputed in a way that preserves the relationship between the predictors and response variables. Furthermore, because many of the lake water quality variables are strongly correlated with each other, their models are expected to be similar. This similarity information can thus be exploited to build better models, especially for the lake water quality variables that have very few observations because they are not sampled frequently. Members of our research team are developing a machine learning approach known as multi-task learning that can simultaneously build regression models of multiple lake water quality variables for a large number of lakes, taking into account both the correlation between the variables and the spatial autocorrelation among the lakes. Because we expect many ecological data sets across broad geographic scales to have similar data gaps and challenges as LAGOS-NE, we think these methods will be extremely valuable for other researchers studying different macroscale questions.

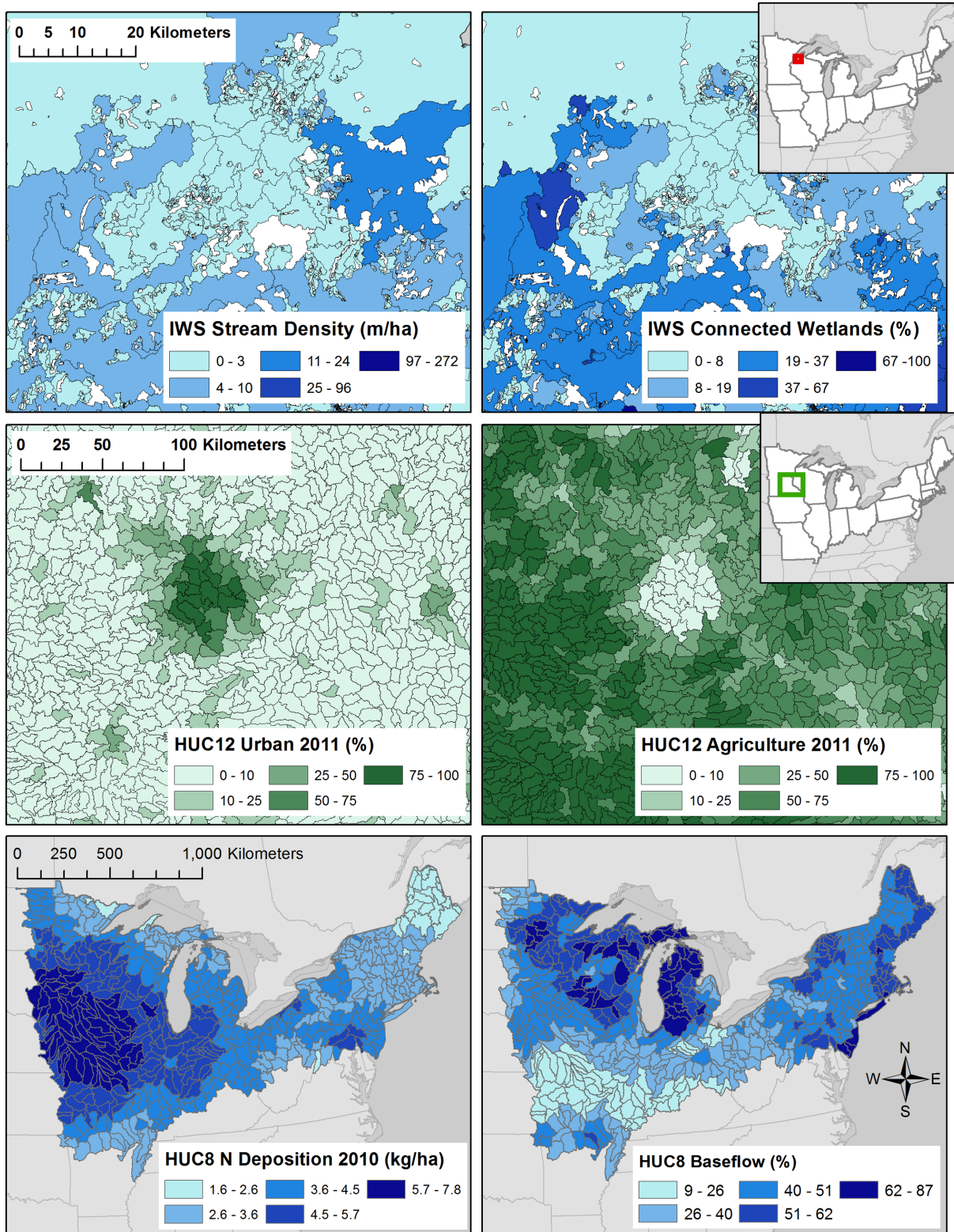


Figure 6: Example ecological context variables by spatial classification in LAGOS-NE. The top 4 panels are zoomed in to selected regions of Minnesota and Wisconsin so that the zone boundaries can be seen. The upper left panel shows stream density in each lake IWS, and the upper right panel shows the percentage of connected wetlands in each lake IWS. The middle left panel shows the 2011 percent urban land use/cover in each hydrologic unit code 12 (HUC12), and the middle right panel shows the 2011 percent agricultural land use/cover in each hydrologic unit code 12 (HUC12). The lower left panel shows the 2010 nitrogen deposition in each HUC8, and the lower right panel shows the average percentage of streamflow that is baseflow in each HUC8.



## Understanding spatial variation in lake nutrients and eutrophication at sub-continental scales

LAGOS-NE allows investigation of spatial variation in lake nutrients and eutrophication at macroscales. For example, members of our team have identified general spatial principles that constrain relationships between ecosystem variables with different spatial structures. In other cases, specific questions regarding spatial patterns have focused on identifying important landscape controls on nutrients and their ratios [38], potential stress induced on phytoplankton communities by high nitrogen levels, and spatial autocorrelation in lake-specific relationships between chlorophyll and nutrients and carbon [39]. In addition, LAGOS-NE contains a wealth of information on a variety of lake ecosystem types. Shallow lakes, in particular, are very abundant across the study area and represent systems that can exhibit hysteresis in response to lake eutrophication. Our team is also investigating the spatial distribution and temporal dynamics of water clarity in shallow lakes of the LAGOS-NE study area.

An important area of research, and one that was a motivating factor for the creation of LAGOS-NE, is understanding the importance of cross-scale interactions (CSIs)—where ecological processes operating at one spatial or temporal scale interact with processes operating at another scale—in lake ecosystems. Because of their importance ecologically and the challenge of quantifying them over large spatial extents, Wagner et al. [40] evaluated the statistical power of large multi-thematic, multi-scaled data sets, such as LAGOS-NE, to detect CSIs. This work not only helped inform the design of large-scale studies aimed at detecting CSIs, but also focused attention on the importance of considering CSI effect sizes and their ecological relevance. To extend this work, members of our team are investigating the importance of both within- and cross-scale interactions in landscape models predicting lake nutrients, and the role that connectivity among freshwaters plays in these interactions. Understanding and predicting nutrients in lakes at macroscales is important to inform estimates of lake contributions to continental and global nutrient cycles. To date, much of this work has been performed on a nutrient-by-nutrient basis, despite knowing that cycles of nitrogen and phosphorus and other key elements are best understood by considering multiple elements in tandem, e.g., in a stoichiometric framework [41] or through analysis of coupled biogeochemical cycles (e.g., [42–44]). Currently, efforts are underway to develop spatial joint nutrient distribution models to evaluate how our understanding of landscape-scale drivers of lake nutrients and predictive performance are improved by considering multiple nutrients simultaneously (multivariate models) compared with traditional univariate approaches that ignore that nutrient cycles can be tightly coupled in freshwaters.

## Understanding temporal and spatial variation in lake eutrophication at sub-continental scales

In addition to the vast spatial data contained in LAGOS-NE, temporal data are available for many water quality variables and some of the ecological context variables (e.g., land use/cover and atmospheric deposition). This is important information within the context of understanding and predicting how lake ecosystems have and will respond to global change, such as changes in climate and land use, and management activities to reduce nutrient inputs to lakes. Because we do not expect responses to such change and actions to be the same everywhere, these questions must be addressed across both space and time. In

particular, recent environmental changes and management efforts have been hypothesized to both improve and degrade water quality in lakes. However, to date, there have been no studies to examine these issues comprehensively across broad scales and to examine which drivers are most strongly related to eutrophication status in lakes. LAGOS-NE is very well suited to answer these types of questions.

For example, nearly 3000 lakes were examined for trends in nutrients and chlorophyll from 1990 to 2013 using LAGOS-NE [45]. Across all lakes, nitrogen has declined, and phosphorus and chlorophyll have not changed. Nitrogen and stoichiometric changes in lakes were related to atmospheric deposition of nitrogen, providing key insight into large-scale nutrient transport and policies such as the Clean Air Act. Using only citizen science data in a subset of the LAGOS-NE database, Lottig et al. [46] showed results that suggested little evidence for major declines or improvements in water quality. In addition, members of our team are examining the relationships between a wide range of climate metrics and water quality in ~11 000 lakes in LAGOS-NE to determine (i) which climate metrics are most related to water quality; (ii) whether physical, chemical, and biological aspects of lakes respond to climate in the same way; and (iii) how the climate–water quality relationship varies across space and regions with different ecological contexts. However, the temporal dynamics of lake ecosystem properties can sometimes be non-linear and exhibit variability across the landscape—largely because of climate and within-lake processes. Our team has developed models for understanding and predicting the often complex temporal patterns observed in water clarity. These studies point to the importance of considering both space and time when trying to understand broad-scale environmental issues in surface waters.

## Using LAGOS-NE for Future Research, Management, and Policy

To facilitate the potential future use of LAGOS-NE, we have thoroughly documented the database and its methods [17], and here we share LAGOS-NE data with the broader research community. In this data paper, we include a wide range of research products, including the water quality and ecological context data, the GIS coverages underlying much of the analyses on freshwaters, and an R package that facilitates use of LAGOS-NE [47]. This package includes functions to retrieve, store, and interact with the LAGOS-NE database, which works across many different operating systems. The package should increase the ease with which users of the database are able to access the data and documentation while maintaining a reproducible workflow.

Key motives for constructing this database included interest in examining lake nutrients and productivity at multiple spatial and temporal scales, fostering broad-scale aquatic ecology and macrosystems research in an open-science platform, and providing new understanding and resources for management and policy-makers. To this end, several team members have made presentations at scientific meetings about the structure and use of LAGOS-NE, and subsets of LAGOS-NE data have been shared with other researchers and stakeholders and agency personnel in advance of this publication. These early uses of LAGOS-NE data by other researchers outside of our team include an investigation of patterns and causes of shifting distribution of a sentinel fish species, developing models to simulate lake temperatures and fish species distributions, and developing a recruitment model for a popular game fish. Results from the latter

2 efforts will inform state-level fisheries management as well as aid in prioritization of lakes for habitat conservation action across a tri-state region.

Much of the research that we and others are conducting with LAGOS-NE has implications for ecosystem management or environmental decision-making. In addition, we have collaborated with boundary organizations and decision-makers. For example, under development is a dashboard of the ecosystem services provided by lakes for use by land managers. In addition, we have helped the state of Michigan determine lake-specific nutrient standards. Our hope is that this database and the associated support tools and documentation serve as a powerful resource and a foundation for future research and decision-making by a broad community of scientists, policy-makers, and natural resource managers. Indeed, our success and experience with database construction and research have inspired us to expand the spatial extent for LAGOS-NE. We have begun to build LAGOS-US, which will include similar data as LAGOS-NE but will be for the continental United States.

### Challenges and Recommendations for Creating Large, Integrated, and Heterogeneous Databases

We found that the largest challenge when creating this database was integrating many small heterogeneous data sets that had few common standards. Although creating such large, integrated data sets using fully automated procedures may happen someday, it appears that we are nowhere near such automation today. Until standards in metadata documentation and robust ontologies are created and widely adopted when creating local or regional data sets, future efforts to integrate these into larger databases will have to rely on close collaborations among domain experts and ecoinformatics professionals, extensive manual interpretation of individual data sets, and funds sufficient to implement these labor-intensive approaches [16]. Nevertheless, it is worth the time and money invested in database integration if the resulting databases support new research, management, policy, public outreach, and education at all levels. We anticipate that LAGOS-NE will serve as a foundation for new data modules that can be used beyond the original intent of LAGOS-NE.

### The economic value of water quality data in an integrated database

This extensive effort was supported by a US National Science Foundation grant that totaled \$2.4 million, along with resources from other projects. Our team ranged in size from 14–20 individuals across the 6 years of the project, with many members compiling and integrating data, authoring metadata, creating new data products, and implementing quality control procedures, resulting in a tremendous number of person-hours. However, when one considers the cost of the data collection for the water quality data in the first place, the expense of this post-processing integration work is not as large as it sounds. Sprague et al. [16] suggest that a single sample (estimated for collecting nutrient or chemistry data from streams) ranged in cost from \$2000 to \$6000 per sample. If we assume similar rates for lake sampling, but lower the cost as some aspects of lake sampling may be cheaper than stream sampling and multiply that cost (estimated as \$1000–\$4000 US) by the total number of records of nutrient or chemical samples in LAGOS-NE ( $n = 589\,909$ ), then the combined estimate to collect the water quality data found in

LAGOS-NE is in the range of \$0.5–2.4 billion US. It cost us between 0.10% and 0.40% of the cost to sample the data in the first place to harmonize these half a million records and to build an ecological context database for them. This relatively small investment in preserving, documenting, and harmonizing these valuable data sets creates the needed infrastructure for new broad-scale research, management, education, and outreach uses.

### Strategies for broad-scale data integration efforts

One challenge is to prioritize research areas and to identify the types of data sets that may benefit from a similar type of integration. State, federal, tribal, and citizen science water quality data sets were an excellent source of quality data for integration and conducting broad-scale research on aquatic systems. There are likely other such data sources that would benefit from being integrated as we have done here. We recommend the following strategies to make the best use of future data integration efforts.

- (1) The database integration effort should be driven by key underlying research questions or goals and grounded in a strong conceptual foundation of the important features to include. In our case, the principles of landscape limnology [12, 18–20] guided the development of LAGOS-NE, which helped us to prioritize geospatial and lake features for inclusion in the database because the addition of any data type or data set cost time and money.
- (2) For databases with more than 1 major data type, it is very helpful to build the database in modular form, each with its own versioning system, specific data integration methods, and quality control procedures. This strategy was not a primary goal at the outset of our project, but it emerged somewhat organically through the life of the project. We now recognize the many benefits that the modularity brings to the database, including making it much easier to be dynamic rather than static by providing a platform for the addition of new data, new types of data, and new modules in the future (such as for biological data or data from high-frequency sensors).
- (3) The entire process should be grounded in an open-science framework. Knowing that the database, design, and methods were to be shared and made usable by future users influenced our decisions throughout the process and made documentation a high priority throughout. Although we are making the full database available now, before this point, we supported open science by publishing subsets of LAGOS-NE data that were used in individual publications (e.g., [48, 49]).
- (4) Creation of LAGOS-NE required a strong focus on team science, and in particular the roles of and incentives for early-career researchers in such efforts. This type of research cannot be conducted in a single-investigator mode, but requires a highly collaborative and effective team-based model (e.g., [50–52]). We explicitly considered strategies for ensuring that early-career team members get credit for their contributions [53], and we recommend providing team members with opportunities for leadership, project management, personnel management, and intellectual growth. For example, they can be part of major decisions and can lead smaller efforts throughout the project, as well as be given power to shape team policies and practices. This integration of early-career researchers into the entire research team and effort will give early-career professionals deep knowledge of the database and procedures, as well as the skills to conduct such work in the future.

- (5) The decision of how to disseminate the database documentation needs to be considered early in the project. For example, database documentation papers are rare, especially in ecology, but are very important. The documentation and procedural approaches for developing this large, integrated, and heterogeneous database had to be disseminated through publication prior to making the database available [17] and prior to publication of research results stemming from LAGOS-NE because methods sections in journal articles are too short to include all the necessary documentation of such methods. Other researchers may be discouraged by the very real consequence that publishing such products takes time and energy investments that may slow down production of research publications. However, such a paper was instrumental in supporting later research articles that used LAGOS-NE. Therefore, we recommend that this (and other) database documentation papers become a more standard type of paper to describe the extensive methods involved and to supplement data papers. Such papers will facilitate the use, extension, and translation of these databases well into the future, as well as foster future research on broad-scale, complex, and societally relevant environmental questions.

### Availability of supporting source code and requirements

Project name: LAGOS-NE

Project home page: <https://github.com/cont-limno/LAGOS>

Operating system(s): e.g., platform independent

Programming language: R

Other requirements: R packages required (with associated versions): dplyr ( $\geq 0.7.0$ ), rappdirs ( $\geq 0.3.1$ ), lazyeval ( $\geq 0.2$ ), purrr ( $\geq 0.2.2.2$ ), magrittr ( $\geq 1.5$ ), sf, curl ( $\geq 2.7.0$ ), stringr ( $\geq 1.2.0$ )

License: GPL

### Availability of supporting data

The data sets supporting the results of this article are available in the Ecological Data Initiative repository, including the following specific components:

- LAGOS-NE-LOCUS v1.01 [54];
- LAGOS-NE-LIMNO v1.087.1 [55];
- LAGOS-NE-GEO v1.05 [56];
- LAGOS-NE-GIS v1.0 [57];
- Snapshots of the R package in the LAGOS GitHub page are also available in the *GigaScience* repository, *GigaDB* [58].

### Additional files

Soranno\_etal.2017\_Additional file 1.8SEP17\_final.docx

Soranno\_etal.2017\_Additional file 2.qaqc-limno.v2.docx

### Abbreviations

CHAG: Climate, Hydrology, Atmospheric deposition of nitrogen and sulfur, and surficial Geology; CONN: connectivity and abundance (lake, stream, and wetland); CSI: cross-scale interactions; DOC: dissolved organic carbon; EML: ecological metadata language; GIS: Geographic Information System; HUC: Hydrologic Unit Code; IQR: interquartile range; IWS: interlake watershed; LAGOS-NE: LAke multi-scaled GeOSpatial and tem-

poral database for the 17 Northeastern and Midwest US states; LULC: land use land cover; MAV: maximum allowable value; NHD: National Hydrography Dataset; SRP: soluble reactive phosphorus; TDN: total dissolved nitrogen; TN: total nitrogen; TP: total phosphorus; US EPA: United States Environmental Protection Agency; USGS: United States Geological Survey; WBD: Watershed Boundary Dataset.

### Competing interests

The authors declare that they have no competing interests.

### Funding

The creation of LAGOS-NE was supported by the National Science Foundation (NSF) MacroSystems Biology Program in the Emerging Frontiers Division of the Biological Sciences Directorate (EF-1065786, EF-1638679, EF-1065649, EF-1065818, EF-1638554) and the United States Department of Agriculture National Institute of Food and Agriculture, Hatch project 176820 to P.A.S. K.E.W. thanks the STRIVE Programme (2011-W-FS-7) from the Environmental Protection Agency, Ireland. S.M.C. thanks the NSF Division of Biological Infrastructure (1401954).

The water quality data that are incorporated into LAGOS-NE were originally funded by the following sources: State of Maine; Michigan Agricultural Experiment Station; Fisheries Division, Michigan Department of Natural Resources; New York State Division of Water Quality; Wisconsin Department of Natural Resources; University of Wisconsin-Madison; State/Trust; Michigan State University Agriculture Experimental Station Disciplinary Research Grant Program; US EPA; US EPA Section 106/319 Grants; Tribal General Fund; US Army Corps of Engineers Federal Lakes Operation and Maintenance Funds; Aquatic Plant Management Society; Aquatic Ecosystem Restoration Foundation; Michigan State University; Michigan State University Department of Fisheries and Wildlife; EPA Star Fellowship to K.S.C. (U-915342-01-0); Andrew W. Mellon Foundation; Federal Aid in Sport Fish Restoration Program (Grant F-69-P, Fish Management in Ohio) administered jointly by the US Fish and Wildlife Service and the Ohio Department of Natural Resources, Division of Wildlife; Iowa Department of Natural Resources (Contract #ESD04HALFasch110155); Minnesota Pollution Control Agency; NSF-Division of Environmental Biology; Ohio Department of Natural Resources Division of Wildlife; University of Rhode Island Watershed Watch; NSF Kellogg Biological Station Long Term Ecological Research (LTER) Program, DEB 1027253; NSF North Temperate Lakes LTER Program, DEB 1440297; Lac du Flambeau Band and Bureau of Indian Affairs; Indiana Department of Environmental Management; Missouri Department of Natural Resources; Clean Water Act Section 16; Michigan Department of Environmental Quality; Massachusetts Water Supply Protection Trust; US EPA Clean Air Markets Division (LTM Network); US EPA Office of Research and Development; New York City Department of Environmental Protection (NYSDEP); City of New York; USGS Water Availability and Use Science Program (WAUSP); US Geological Survey; New York State Energy Research and Development Authority; National Institute of Food and Agriculture, US Department of Agriculture, Hatch Grant 1003732; the New York State Department of Environmental Conservation; Lake Sunapee Protective Association; National Oceanic and Atmospheric Administration; Gull Lake Quality Organization; Clean Michigan Initiative; NSF grant DEB-1455461.

## Author contributions

Data for the database were contributed by L.C.B., M.B., K.E.B., M.G.B., M.T.B., S.R.C., J.W.C., K.S.C., M.C., J.D.C., J.A.D., J.D., C.T.F., C.S.F., M.J.G., L.T.G., J.D.F., S.K.H., P.C.H., E.H., C.H., J.R.J., K.J.H., L.L.J., W.W.J., J.R.J., C.M.K., S.A.K., B.L., J.A.L., Y.L., N.R.L., J.A.L., L.J.M., W.H.M., K.E.B.M., B.P.N., S.J.N., M.L.P., D.C.P., A.I.P., D.M.P., P.O.R., D.O.R., K.M.R., L.G.R., O.S., N.J.S., P.A.S., N.R.S., E.H.S., J.L.S., J.M.T., T.P.T., M.V., G.W., K.C.W., K.E.W., J.D.W., and M.K.W. The idea to create the database was conceived by P.A.S. and K.S.C. P.A.S. coordinated the different activities across team members to build LAGOS-NE. The database was designed by E.G.B., P.N.T., C.G., and P.A.S. and created and managed by E.G.B. The following authored metadata for the individual water quality data sets using information provided by the data providers: M.T.B., C.K.B., K.S.C., S.M.C., C.E.F., C.T.F., E.N.H., N.R.L., S.K.O., N.K.S., P.A.S., E.H.S., and K.E.W. C.E.F. prepared the integrated LAGOS-NE metadata and developed the protocols for authoring the EML metadata, and C.E.F. and C.K.B. created EML metadata for the 87 water quality data sets. S.K.O. wrote the final variables' definitions for the integrated metadata. C.G. helped to prepare the needed metadata and documentation for loading the data in the data repository. Code for importing the data sets into the database was written by E.G.B., S.T.C., N.R.L., and S.Y. N.J.S. and S.B.S. performed geospatial analyses and created the LAGOS-GIS Toolbox. The conceptual foundation for measuring freshwater connectivity was led by C.E.F. S.B.S. developed the methods to delineate lake watersheds. The quality control methods development and analysis on LAGOS-NE<sub>LIMNO</sub> were conducted by N.R.L.; the quality control of LAGOS-NE<sub>GIS</sub> was led by C.E.S. and S.M.C. and conducted by C.E.S., S.M.C., C.E.F., N.K.S., and K.E.W. The quality control of LAGOS-NE<sub>LOCUS</sub> was conducted by E.G.B. Many authors who were part of the database integration team wrote the technical documentation; J.F.L. served as editor of these technical documents. Tables and figures were prepared by S.M.C., K.B.S.K., J.F.L., N.R.L., A.C.P., N.K.S., and P.A.S. and edited by many of the contributing authors. S.K.O. and J.J.S. wrote the LAGOS-NE R package. N.J.S. prepared the GIS data and their corresponding metadata. P.A.S. coordinated the writing of the manuscript, and major parts of the manuscript were written by P.A.S., K.S.C., S.M.C., J.F.L., N.R.L., S.K.O., J.J.S., E.H.S., P.N.T., T.W., and S.Y. After the lead author, authors are listed alphabetically.

## Acknowledgments

We thank the contributions over the past several decades of many hundreds to thousands of governmental, tribal, and citizen scientists whose efforts from lake sampling to water quality analysis to data set compilation enabled LAGOS-NE to become a reality and a resource for the future. Specifically, we dedicate this paper to the memory of Jody Connor, whose 3 decades of innovative and science-driven lake management while working for the New Hampshire Department of Environmental Services generated meaningful contributions to the protection and restoration of lake quality and to lasting data legacies such as volunteer monitoring in the state. This is Great Lakes Environmental Research Laboratory contribution number 1863. Any use of trade, firm, or product names is for descriptive purposes only and does not imply endorsement by the US Government.

## References

- Carpenter SR, Caraco NF, Correll DL et al. Nonpoint pollution of surface waters with phosphorus and nitrogen. *Ecol Appl* 1998;**8**(3):559–68.
- Jaworski NA, Howarth RW, Hetling LJ. Atmospheric deposition of nitrogen oxides onto the landscape contributes to coastal eutrophication in the Northeast United States. *Environ Sci Technol* 1997;**31**(7):1995–2004.
- Bennett EM, Carpenter SR, Caraco NF. Human impact on erodible phosphorus and eutrophication: a global perspective. *Bioscience* 2001;**51**(3):227–34.
- Schindler DW. Recent advances in the understanding and management of eutrophication. *Limnol Oceanogr* 2006;**51**(1part2):356–63.
- Taranu ZE, Gregory-Eaves I. Quantifying relationships among phosphorus, agriculture, and lake depth at an inter-regional scale. *Ecosystems* 2008;**11**(5):715–25.
- Filstrup CT, Wagner T, Soranno PA et al. Regional variability among nonlinear chlorophyll-phosphorus relationships in lakes. *Limnol Oceanogr* 2014;**59**(5):1691–703.
- Mccrackin ML, Jones HP, Jones PC et al. Recovery of lakes and coastal marine ecosystems from eutrophication: a global meta-analysis. *Limnol Oceanogr* 2017;**62**(2):507–18.
- Paerl HW, Otten TG, Joyner AR. Moving towards adaptive management of cyanotoxin-impaired water bodies. *Microb Biotechnol* 2016;**9**(5):641–51.
- Schindler DW, Carpenter SR, Chapra SC et al. Reducing phosphorus to curb lake eutrophication is a success. *Environ Sci Technol* 2016;**50**(17):8923–9.
- Emi Fergus C, Soranno PA, Cheruvilil KS et al. Multiscale landscape and wetland drivers of lake total phosphorus and water color. *Limnol Oceanogr* 2011;**56**(6):2127–46.
- Soranno PA, Cheruvilil KS, Bissell EG et al. Cross-scale interactions: quantifying multi-scaled cause-effect relationships in macrosystems. *Front Ecol Environ* 2014;**12**(1):65–73.
- Read EK, Patil VP, Oliver SK et al. The importance of lake-specific characteristics for water quality across the continental United States. *Ecol Appl* 2015;**35**(4):943–55.
- Smith VH, Dodds WK, Havens KE et al. Comment: cultural eutrophication of natural lakes in the United States is real and widespread. *Limnol Oceanogr* 2014;**59**(6):2217–25.
- McDonald CP, Lottig NR, Stoddard JL et al. Comment on Bachmann (2013): a non-representative sample cannot describe the extent of cultural eutrophication of natural lakes in the United States. *Limnol Oceanogr* 2014;**59**:2226–30.
- Stoddard JL, Van Sickle J, Herlihy AT et al. Continental-scale increase in lake and stream phosphorus: are oligotrophic systems disappearing in the United States? *Environ Sci Technol* 2016;**50**(7):3409–15.
- Sprague LA, Oelsner GP, Argue DM. Challenges with secondary use of multi-source water-quality data in the United States. *Water Res* 2017;**100**:252–61.
- Soranno PA, Bissell EG, Cheruvilil KS et al. Building a multi-scaled geospatial temporal ecology database from disparate data sources: fostering open science and data reuse. *Giga-science* 2015;**4**(1):28.
- Magnuson JJ, Kratz TK. Lakes in the landscape: approaches to regional limnology. *Int Assoc Theoret Appl Limnol* 2000;**27**:74–87.
- Wiens JA. Riverine landscapes: taking landscape ecology into the water. *Freshwater Biol* 2002;**47**(4):501–15.
- Soranno PA, Cheruvilil KS, Webster KE et al. Using landscape limnology to classify freshwater ecosystems for multi-ecosystem management and conservation. *Bioscience* 2010;**60**(6):440–54.
- United States Geological Survey national hydrography dataset. Version 9.3. <http://nhd.usgs.gov>. Accessed 4 June 2015.

22. Zhang T, Soranno PA, Cheruvilil KS et al. Evaluating the effects of upstream lakes and wetlands on lake phosphorus concentrations using a spatially-explicit model. *Landscape Ecol* 2012;27(7):1015–30.
23. United States Geological Survey watershed boundary dataset. <https://nhd.usgs.gov/wbd.html>. Accessed 2013.
24. National elevation dataset. <http://ned.usgs.gov/>. Accessed 11 March 2013.
25. US Environmental Protection Agency: national lakes assessment fact sheet. 2010. [http://water.epa.gov/type/lakes/upload/nla\\_survey\\_fact\\_sheet.pdf](http://water.epa.gov/type/lakes/upload/nla_survey_fact_sheet.pdf). Accessed 4 June 2015.
26. US Environmental Protection Agency: national lakes assessment 2012: a fact sheet for communities. 2012. <http://water.epa.gov/type/lakes/assessmonitor/lakessurvey/upload/NLA-2012-Fact-Sheet-for-Communities.pdf>. Accessed 4 June 2015.
27. Environmental Protection Agency: National Lake Survey of 2012. 2012. <https://www.epa.gov/national-aquatic-resource-surveys/nla>. Accessed 4 June 2015.
28. Oliver SK, Soranno PA, Fergus CE et al. Prediction of lake depth across a 17-state region in the United States. *Inland Waters* 2016;6:314–24.
29. Fergus CE, Lapierre J, Oliver SK et al. The freshwater landscape: lake, wetland, and stream abundance and connectivity at macroscales. *Ecosphere* 2017;8(8):e01911.
30. Wagner T, Soranno PA, Cheruvilil KS et al. Quantifying sample biases of inland lake sampling programs in relation to lake surface area and land use/cover. *Environ Monit Assess* 2008;141(1-3):131–47.
31. Strock KE, Saros JE, Nelson SJ et al. Extreme weather years drive episodic changes in lake chemistry: implications for recovery from sulfate deposition and long-term trends in dissolved organic carbon. *Biogeochemistry* 2016;127(2-3):353–65.
32. Seaber PR, Kapinos FP, Knapp GL. Hydrologic unit maps: US Geological Survey water-supply paper 2294. 1987. <http://water.usgs.gov/GIS/huc.html>. Accessed 11 March 2013.
33. Smith NJ, Soranno PA, Stopyak S. LAGOS-NE GIS Toolbox. GitHub 2014. <https://soranno.github.io/LAGOS-GIS-Toolbox/>. Accessed 1 December 2016.
34. Yuan S, Tan PN, Cheruvilil KS et al. Constrained spectral clustering for regionalization: exploring the trade-off between spatial contiguity and landscape homogeneity. *Data Sci Adv Analyt* 2015; doi:10.1109/DSAA.2015.7344878.
35. Cheruvilil KS, Yuan S, Webster KE et al. Creating multi-themed ecological regions for macroscale ecology: testing a flexible, repeatable, and accessible clustering method. *Ecol Evol* 2017;7(9):3046–58.
36. Helsel DR. *Statistics for Censored Environmental Data Using Minitab and R*, 2nd edition. New York: John Wiley and Sons; 2012.
37. Yuan S, Tan PN, Cheruvilil KC et al. Hash-based feature learning for incomplete continuous-valued data. In: *Proceedings of the 2017 SIAM International Conference on Data Mining*. Houston, TX, 2017. p. 678–686. Philadelphia, PA: Society for Industrial and Applied Mathematics, 2017.
38. Collins SM, Oliver SK, Lapierre J et al. Lake nutrient stoichiometry is less predictable than nutrient concentrations at regional and sub-continental scales. *Ecol Appl* 2017;27(5):1529–40.
39. Fergus CE, Finley AO, Soranno PA et al. Spatial variation in nutrient and water color effects on lake chlorophyll at macroscales. *PLoS One* 2016;11(10):e0164592.
40. Wagner T, Fergus CE, Stow CA et al. The statistical power to detect cross-scale interactions at macroscales. *Ecosphere* 2016;7:e01417.
41. Sterner RW, Elser JJ. *The Biology of Elements from Molecules to the Biosphere*. Princeton, NJ: Princeton University Press; 2002.
42. Rastetter EB. Modeling coupled biogeochemical cycles. *Front Ecol Environ* 2011;9:68–73.
43. Finzi AC, Austin AT, Cleland EE et al. Responses and feedbacks of coupled biogeochemical cycles to climate change: examples from terrestrial ecosystems. *Front Ecol Environ* 2011;9:61–67.
44. Finlay JC, Small GE, Sterner RW. Human influences on nitrogen removal in lakes. *Science* 2013;342:247–50.
45. Oliver SK, Collins SM, Soranno PA et al. Unexpected stasis in a changing world: lake nutrient and chlorophyll trends since 1990. *Glob Change Biol* 2017; doi:10.1111/gcb.13810.
46. Lottig NR, Wagner T, Norton Henry E et al. Long-term citizen-collected data reveal geographical patterns and temporal trends in lake water clarity. *PLoS One* 2014; doi: 10.1371/journal.pone.0095769.
47. Stachelek J, Oliver SK. LAGOS: R interface to the LAke multi-scaled GeOSpatial & temporal database. R package version 1.087.1. GitHub 2017. <https://github.com/cont-limno/LAGOS>. Accessed 1 September 2017.
48. Oliver SK, Soranno PA, Fergus CE et al. LAGOS – predicted and observed maximum depth values for lakes in a 17-state region of the U.S. *Long Term Ecological Research Network* 2015; doi:10.6073/pasta/f00a245fd9461529b8cd9d992d7e3a2f. Accessed 1 September 2017.
49. Fergus CE, Finley AO, Soranno PA et al. Spatial variation in nutrient and water color effects on lake chlorophyll at macroscales. *Long-Term Ecological Research Network Data Portal* 2016; doi: 10.6073/pasta/0ebd2e4c0705706b77b359955bff44e1. Accessed 1 September 2017.
50. Cheruvilil KS, Soranno PA, Weathers KC et al. Creating and maintaining high-performing collaborative research teams: the importance of diversity and interpersonal skills. *Front Ecol Environ* 2014;12:31–38.
51. Weathers KC, Hanson PC, Arzberger P et al. The Global Lake Ecological Observatory Network (GLEON): the evolution of grassroots network science. *Bull Limnol Oceanogr* 2013;22:71–73.
52. Hanson PC, Weathers KC, Kratz TK. Networked lake science: how the Global Lake Ecological Observatory (GLEON) works to understand, predict, and communicate lake ecosystem response to global change. *Inland Waters* 2016; doi: 10.5268/IW-6.4.904.
53. Goring SJ, Weathers KC, Dodds WK et al. Improving the culture of interdisciplinary collaboration in ecology by expanding measures of success. *Front Ecol Environ* 2014;14: 39–47.
54. Soranno P, Cheruvilil K. LAGOS-NE-LOCUS v1.01: a module for LAGOS-NE, a multi-scaled geospatial and temporal database of lake ecological context and water quality for thousands of U.S. Lakes: 1925–2013. *Environmental Data Initiative*. 2017; <http://dx.doi.org/doi:10.6073/pasta/0c23a789232ab4f92107e26f70a7d8ef>. Accessed 10 October 2017.
55. Soranno P, Cheruvilil K. LAGOS-NE-LIMNO v1.087.1: a module for LAGOS-NE, a multi-scaled geospatial and temporal database of lake ecological context and water quality for thousands of U.S. Lakes: 1925–2013. *Environmental*

- Data Initiative. 2017; <http://dx.doi.org/10.6073/pasta/56cc5f1f753d48edfea170a5401dd6df>. Accessed 10 October 2017.
56. Soranno P, Cheruvilil K. LAGOS-NE-GEO v1.05: a module for LAGOS-NE, a multi-scaled geospatial and temporal database of lake ecological context and water quality for thousands of U.S. Lakes: 1925–2013. Environmental Data Initiative. 2017; <http://dx.doi.org/doi:10.6073/pasta/16f4bdaa9607c845c0b261a580730a7a>. Accessed 10 October 2017.
57. Soranno P, Cheruvilil K. LAGOS-NE-GIS v1.0: a module for LAGOS-NE, a multi-scaled geospatial and temporal database of lake ecological context and water quality for thousands of U.S. Lakes: 2013–1925. Environmental Data Initiative. 2017; <http://dx.doi.org/doi:10.6073/pasta/8674fd113c0089c0fa174ee4eaf3f376>. Accessed 10 October 2017.
58. Soranno PA, Bacon LC, Beauchene M et al. Supporting data for “LAGOS-NE: a multi-scaled geospatial and temporal database of lake ecological context and water quality for thousands of US lakes.” GigaScience Database 2017. <http://dx.doi.org/10.5524/100350>. Accessed 1 October 2017.