

GigaScience

LAGOS-NE: A multi-scaled geospatial and temporal database of lake ecological context and water quality for thousands of U.S. lakes

--Manuscript Draft--

Manuscript Number:	GIGA-D-17-00112R2	
Full Title:	LAGOS-NE: A multi-scaled geospatial and temporal database of lake ecological context and water quality for thousands of U.S. lakes	
Article Type:	Data Note	
Funding Information:	State of Maine	Lake Ecologist Linda C Bacon
	U.S. Environmental Protection Agency (106/319)	Not applicable
	U.S. Army Corps of Engineers (US)	Limnologist Marvin G Boyer
	Michigan Agricultural Experiment Station Disciplinary Research Grant Program	Associate Professor Mary T Bremigan
	Michigan Department of Natural Resources	Professor Kendra S Cheruvellil
	National Science Foundation (1027253)	Professor Stephen K Hamilton
	New York State Division of Water Quality	Chief Scott A Kishbaugh
	Wisconsin Department of Natural Resources	Professor Stephen R Carpenter
	University of Wisconsin-Madison	Professor Stephen R Carpenter
	National Science Foundation (1455461)	Professor Stephen R Carpenter
	Tribal General Fund	Water Regulatory and Restoration Specialist Karen E Bednar
	Aquatic Plant Management Society	Associate Professor Mary T Bremigan
	Aquatic Ecosystem Restoration Foundation	Associate Professor Mary T Bremigan
	Michigan State University	Associate Professor Mary T Bremigan
	U.S. Environmental Protection Agency (U-915958801-2)	Professor Kendra S Cheruvellil
	Andrew W. Mellon Foundation	Professor Stephen R Carpenter
	U.S. Fish and Wildlife Service (F-69-R)	Dr. Joseph D Conroy
	Iowa Department of Natural Resources (ESD04HALFasch110155)	Research Associate Christopher T Filstrup
	Minnesota Pollution Control Agency	Professor John A Downing
	National Science Foundation Macrosystems Biology Program (EF-1065786, EF-1638679)	Professor Patricia Ann Soranno
	Ohio Department of Natural Resources Division of Wildlife	Professor Maria J Gonzalez
	University of Rhode Island Watershed Watch	Program Director Linda T Green
	National Science Foundation North Temperate Lakes LTER Program (1440297)	Research Professor Paul C Hanson
	Lac du Flambeau Band and Bureau of Indian Affairs	Ecologist Celeste Hockings
	Indiana Department of Environmental Management	Professor Emeritus William W Jones
	Missouri Department of Natural Resources	Curators' Professor Emeritus John R Jones

U.S. Environmental Protection Agency (US) Clean Water Act Section 16	Water Quality Specialist Caroline M Keson
Michigan Department of Environmental Quality	Academic Specialist Jo A Latimore
Massachusetts Water Supply Protection Trust	Environmental Analyst Yuehlin Lee
U.S. Environmental Protection Agency (US) Clean Air Markets Division (LTM Network)	Biologist Clara S Funk
New York City Department of Environmental Protection	Research Scientist Karen E.B Moore
U.S. Geological Survey (US) Water Availability and Use Science Program (WAUSP)	Research Hydrologist Brian P Neff
U.S. Geological Survey	Research Hydrologist Donald O Rosenberry
New York State Energy Research and Development Authority	MS Karen M Roy
National Institute of Food and Agriculture (1003732)	Professor Lars G Rudstam
New York State Department of Environmental Conservation	Senior Research Associate James R Jackson
Lake Sunapee Protective Association	G. Evelyn Hutchinson Chair of Ecology Kathleen C Weathers
National Oceanic and Atmospheric Administration	Senior Scientist Craig A Stow
Gull Lake Quality Organization	Senior Scientist Craig A Stow
Clean Michigan Initiative	Aquatic Biologist Marcy K Wilmes
City of New York	Chief Lorraine L Janus
National Science Foundation (US) DBI (1401954)	Postdoctoral Fellow Sarah M Collins
National Institute of Food and Agriculture (176820)	Professor Patricia Ann Soranno
U.S. Environmental Protection Agency (US) Office of Research and Development	Not applicable

Abstract:

Background

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 datasets across broad regions and across time into an integrated database. LAGOS-NE accomplishes this goal for lakes in the northeastern-most 17 U.S. states.

Findings

LAGOS-NE contains data for 51,101 lakes and reservoirs larger than 4 ha in 17 lake-rich U.S. states. The database includes three 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 three decades for approximately 2,600-12,000 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 datasets from federal, state, tribal, and non-profit agencies, university researchers, and citizen scientists.

Conclusions

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.
Corresponding Author:	Patricia Ann Soranno, Ph.D Michigan State University East Lansing, Michigan UNITED STATES
Corresponding Author Secondary Information:	
Corresponding Author's Institution:	Michigan State University
Corresponding Author's Secondary Institution:	
First Author:	Patricia Ann Soranno, Ph.D
First Author Secondary Information:	
Order of Authors:	Patricia Ann Soranno, Ph.D
	Linda C Bacon
	Michael Beauchene
	Karen E Bednar
	Edward G Bissell
	Claire K Boudreau
	Marvin G Boyer
	Mary T Bremigan
	Stephen R Carpenter
	Jamie W Carr
	Kendra S Cheruvellil
	Samuel T Christel
	Matt Claucherty
	Sarah M Collins
	Joseph D Conroy
	Jed Dukett
	John A Downing
	C. Emi Fergus
	Christopher T Filstrup
	Clara S Funk
	Maria J Gonzalez
	Linda T Green
	Corinna Gries
	John D Halfman
	Stephen K Hamilton
	Paul C Hanson
	Emily N Henry
	Elizabeth M Herron

Celeste Hockings
James R Jackson
Kari Jacobson-Hedin
Lorraine L Janus
William W Jones
John R Jones
Caroline M Keson
Katelyn B.S King
Scott A Kishbaugh
Jean-Francois Lapierre
Barbara F Lathrop
Jo A Latimore
Yuehlin Lee
Noah R Lottig
Jason A Lynch
Leslie J Matthews
Bill H McDowell
Karen E.B Moore
Brian P Neff
Sarah J Nelson
Samantha K Oliver
Michael L Pace
Don C Pierson
Autumn C Poisson
Amina I Pollard
David M Post
Paul O Reyes
Donald O Rosenberry
Karen M Roy
Lars G Rudstam
Orlando Sarnelle
Nancy J Schuldt
Caren E Scott
Nicholas K Skaff
Nicole J Smith
Nick R Spinelli
Joseph J Stachelek
Emily H Stanley
John L Stoddard
Scott B Stopyak
Craig A Stow

	Jason M Tallant
	Pang-Ning Tan
	Anthony P Thorpe
	Ty Wagner
	Michael J Vanni
	Gretchen Watkins
	Kathleen C Weathers
	Katherine E Webster
	Jeffrey D White
	Marcy K Wilmes
	Shuai Yuan
Order of Authors Secondary Information:	
Response to Reviewers:	<p>All suggested changes have been made. Please note that some author affiliations have changed, which I updated in the word file, however, i did not see where to change/add author affiliations to the author-boxes in the online form.</p> <p>I have also made some minor changes to text and figure/table legends based on an internal USGS review that i was required to do because I have numerous USGS-affiliated co-authors. If you would like to see a track-changes version of the word doc that shows those changes, I can provide that. However, they were mostly minor editorial changes.</p>
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
Experimental design and statistics	Yes
<p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	
Resources	Yes
<p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p>	

<p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>

1
2
3
4 1 *For submission as a Data Note to GigaScience*

5 2
6 3
7 4
8 5
9 6
10 7
11 7 **LAGOS-NE: A multi-scaled geospatial and temporal database of lake ecological**
12 8 **context and water quality for thousands of U.S. lakes**
13 9
14 9

15 10
16 11
17 12 **Authors:**

18 13 Patricia A. Soranno¹, Linda C. Bacon², Michael Beauchene³, Karen E. Bednar⁴, Edward G. Bissell¹, Claire K.
19 14 Boudreau¹, Marvin G. Boyer⁵, Mary T. Bremigan¹, Stephen R. Carpenter⁶, Jamie W. Carr⁷, Kendra S.
20 15 Cheruvelil¹, Samuel T. Christel⁶, Matt Claucherty⁸, Sarah M. Collins⁶, Joseph D. Conroy⁹, John A. Downing¹⁰,
21 16 Jed Dukett¹¹, C. Emi Fergus¹², Christopher T. Filstrup¹⁰, Clara Funk¹³, Maria J. Gonzalez¹⁴, Linda T. Green¹⁵,
22 17 Corinna Gries⁶, John D. Halfman¹⁶, Stephen K. Hamilton¹⁷, Paul C. Hanson⁶, Emily N. Henry¹⁸, Elizabeth M.
23 18 Herron¹⁹, Celeste Hockings²⁰, James R. Jackson²¹, Kari Jacobson-Hedin²², Lorraine L. Janus²³, William W.
24 19 Jones²⁴, John R. Jones²⁵, Caroline M. Keson²⁶, Katelyn B.S. King¹, Scott A. Kishbaugh²⁷, Jean-Francois
25 20 Lapierre²⁸, Barbara Lathrop²⁹, Jo A. Latimore¹, Yuehlin Lee³⁰, Noah R. Lottig³¹, Jason A. Lynch¹³, Leslie J.
26 21 Matthews³³, William H. McDowell³⁴, Karen E.B. Moore³⁵, Brian P. Neff³⁶, Sarah J. Nelson³⁷, Samantha K.
27 22 Oliver⁶, Michael L. Pace³⁸, Donald C. Pierson³⁹, Autumn C. Poisson¹, Amina I. Pollard⁴⁰, David M. Post⁴¹, Paul
28 23 O. Reyes³⁰, Donald O. Rosenberry⁴², Karen M. Roy⁴³, Lars G. Rudstam⁴⁴, Orlando Sarnelle¹, Nancy J.
29 24 Schuldt⁴⁵, Caren E. Scott⁴⁶, Nicholas K. Skaff¹, Nicole J. Smith¹, Nick R. Spinelli⁴⁷, Joseph J. Stachelek¹, Emily
30 25 H. Stanley⁶, John L. Stoddard⁴⁸, Scott B. Stopyak⁴⁹, Craig A. Stow⁵⁰, Jason M. Tallant⁵¹, Pang-Ning Tan⁵²,
31 26 Anthony P. Thorpe²⁵, Michael J. Vanni⁵³, Tyler Wagner⁵⁴, Gretchen Watkins⁴, Kathleen C. Weathers⁵⁶,
32 27 Katherine E. Webster⁵⁷, Jeffrey D. White⁵⁸, Marcy K. Wilmes⁵⁹, Shuai Yuan⁵²
33 28
34 29

35 26
36 27
37 28
38 29 **ORCID:**

39 30 Patricia Soranno, 0000-0003-1668-9271
40 31 Kendra Cheruvelil, 0000-0003-1880-2880
41 32 Sarah Collins, 0000-0001-5503-7386
42 33 C.Emi Fergus, 0000-0001-8722-9788
43 34 Jean Francois Lapierre, 0000-0001-5862-7955
44 35 Nicholas Skaff, 000-0002-5929-3966
45 36 Pang-Ning Tan, 0000-0003-3205-0339
46 37
47 38
48 39

49 38
50 39 ¹Department of Fisheries and Wildlife, Michigan State University, East Lansing, MI 48824, USA

51 40 ²Department of Environmental Protection, State of Maine, Augusta, ME 04330, USA

52 41 ³Department of Energy and Environmental Protection, State of Connecticut, Hartford, CT 06106, USA

53 42 ⁴Water Resources Program, Lac du Flambeau Tribal Natural Resources, Lac du Flambeau, WI, USA

54 43 ⁵Environmental Planning, US Army Corps of Engineers, Kansas City, MO 64106, USA

55 44 ⁶Center for Limnology, University of Wisconsin Madison, Madison, WI 53706 USA

56 45 ⁷Office of Watershed Management, Massachusetts Department of Conservation and Recreation, West Boylston, MA 10583, USA

57 46 ⁸Watershed Protection, Tipp of the Mitt Watershed Council, Petoskey, MI 49770, USA

58 47 ⁹Division of Wildlife, Inland Fisheries Research Unit, Ohio Department of Natural Resources, Hebron, OH 43025, USA

59 48 ¹⁰Large Lakes Observatory, University of Minnesota, Duluth, MN 55812 USA
60
61
62
63
64
65

1
2
3
4 49 ¹¹Adirondack Lake Survey Corporation, Ray Brook, NY 12977 USA
5 50 ¹²National Research Council, US Environmental Protection Agency, Corvallis, OR 97333, USA
6 51 ¹³Office of Air and Radiation, US Environmental Protection Agency, Washington, DC 20460, USA
7 52 ¹⁴Department of Biology, Miami University, Oxford, OH 45056, USA
8 53 ¹⁵Natural Resource Science, University of Rhode Island, Kingston, RI 02892 USA
9 54 ¹⁶Geoscience, Hobart & William Smith Colleges, Geneva, NY 14456 USA
10 55 ¹⁷Kellogg Biological Station, Michigan State University, Hickory Corners, MI 49060, USA
11 56 ¹⁸Outreach and Engagement, Oregon State University, Corvallis, OR 97331, USA
12 57 ¹⁹Watershed Watch, University of Rhode Island, Kingston, RI 02881, USA
13 58 ²⁰Natural Resource Department, Lac du Flambeau Band of Lake Superior Chippewa Indians, Lac du Flambeau, WI 54538, USA
14 59 ²¹Department of Natural Resources, Cornell University, Bridgeport, NY, USA
15 60 ²²Office of Water Protection, Fond du Lac Reservation, Cloquet, MN 55720 USA
16 61 ²³Bureau of Water Supply, New York City Department of Environmental Protection, Valhalla, NY 10560, USA
17 62 ²⁴School of Public and Environmental Affairs, Indiana University, Bloomington, IN 47408, USA
18 63 ²⁵School of Natural Resources, University of Missouri, Columbia, MO, USA
19 64 ²⁶Natural Resource Department, Little Traverse Bay Bands of Odawa Indians, Harbor Springs, MI 49740, USA
20 65 ²⁷Division of Water, New York State Department of Environmental Conservation, Albany, NY 12233, USA
21 66 ²⁸Department of Biological Science, University of Montreal, Montreal Quebec, Canada, H3C 3J7
22 67 ²⁹Pennsylvania Department of Environmental Protection, State of Pennsylvania, Harrisburg, PA 17101 USA
23 68 ³⁰Office of Watershed Management, Massachusetts Department of Conservation and Recreation, Belchertown, MA 01007, USA
24 69 ³¹Trout Lake Research Station, University of Wisconsin, Boulder Junction, WI 54512, USA
25 70 ³³Lakes and Ponds Program, Vermont Department of Environmental Conservation, Montpelier, VT 05620, USA
26 71 ³⁴Natural Resources and the Environment, University of New Hampshire, Durham, NH 03824, USA
27 72 ³⁵Water Quality Science and Research, New York City Department of Environmental Protection, Kingston, NY 12401, USA
28 73 ³⁶National Research Program, USGS, Denver CO 80225, USA
29 74 ³⁷School of Forest Resources, University of Maine, Orono, ME, USA
30 75 ³⁸Department of Environmental Science, University of Virginia, Charlottesville, VA 22904, USA
31 76 ³⁹Department of Ecology and Genetics, Uppsala University, Uppsala, Sweden
32 77 ⁴⁰Office of Water, US EPA, Washington, DC 20460, USA
33 78 ⁴¹Ecology and Evolutionary Biology, Yale University, Connecticut 06511, USA
34 79 ⁴²National Research Program, USGS, Denver, CO 80225, USA
35 80 ⁴³Division of Air Resources, New York State Department of Environmental Conservation, Ray Brook, NY 12977, USA
36 81 ⁴⁴Department of Natural Resources, Cornell University, Ithaca, NY 14850, USA
37 82 ⁴⁵Environmental Program, Fond du Lac Band of Lake Superior Chippewa Indians, Cloquet, MN 55720, USA
38 83 ⁴⁶Aquatic Science, NEON, Boulder, CO 80301, USA
39 84 ⁴⁷Watershed Management, Lake Wallenpaupack Watershed Management District, Hawley, PA, USA
40 85 ⁴⁸Western Ecology Division, Office of Research and Development, US EPA, Corvallis, OR 97333, USA
41 86 ⁴⁹Technology Services, Eaton County, Charlotte, MI, USA
42 87 ⁵⁰Great Lakes Environmental Research Lab, NOAA, Ann Arbor, MI 47176, USA
43 88 ⁵¹Biological Station, University of Michigan, Pellston, MI 49769, USA
44 89 ⁵²Computer Science and Engineering, Michigan State University, East Lansing, MI 48824, USA
45 90 ⁵³Department of Zoology, Miami University, Oxford, OH 45056 USA
46 91 ⁵⁴ Pennsylvania Cooperative Fish and Wildlife Research Unit, USGS, 402 Forest Resources Building, University Park, PA 16802,
47 92 USA.
48 93 ⁵⁶ Cary Institute of Ecosystem Studies, Millbrook, NY, USA
49 94 ⁵⁷School of Natural Sciences, Trinity College, Dublin, Ireland
50 95 ⁵⁸Biology Department, Framingham State University, Framingham, MA 01702, USA
51 96 ⁵⁹Department of Environmental Quality, State of Michigan, Lansing, MI 48909, USA
52 97
53 98
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

ABSTRACT

Background

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 datasets across broad regions and across time into an integrated database. LAGOS-NE accomplishes this goal for lakes in the northeastern-most 17 U.S. states.

Findings

LAGOS-NE contains data for 51,101 lakes and reservoirs larger than 4 ha in 17 lake-rich U.S. states. The database includes three 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 three decades for approximately 2,600-12,000 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 datasets from federal, state, tribal, and non-profit agencies, university researchers, and citizen scientists.

Conclusions

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

1. 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 [2,3,4]; 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., [5,6,7]). 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., [10,11,12]). 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 concentrations) and measures of lake ecological context across regions, the continent, and the globe (e.g., 13,14,15).

In practice, measures of water quality are often collected from a relatively small number of lakes within individual regions. In the U.S., large investments have been made in water quality monitoring by federal, state, local, and tribal governments; and, many, but not all, of the datasets 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 over 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 and northeastern U.S. (Figure 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 ≥ 4 ha in the study area (1,800,000 km²).

We had three related motivations for developing this database: (1) 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); (2) to build the capacity to apply this scientific understanding to environmental management and policy of inland waters; and, (3) 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.

Figure 1. Map of the study extent of LAGOS-NE. Map includes 17 states in the upper midwest and northeastern U.S. outlined in white and 51,101 lakes ≥ 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 [17].

1
2
3
4 195
5 196 LAGOS-NE comprises three data modules that, although integrated in the same database, were
6 197 derived using different data sources and data integration methods, and thus must be version-controlled
7 198 separately. LAGOS-NE_{LOCUS} v1.01 includes lake location and physical characteristics based on an
8 199 existing national-scale database of lake and streams in the U.S. for all lakes. LAGOS-NE_{GEO} v1.05
9 200 includes measures of land, water, and air (ecological context) obtained from existing national scale GIS
10 201 (geographic information system) datasets and measured in multiple zones (delineated by different spatial
11 202 classifications) around all lakes. This module also contains some temporal data for climate, land
12 203 use/cover, and atmospheric deposition variables. LAGOS-NE_{LIMNO} v1.087.1 includes in-situ
13 204 measurements of lake water quality for a subset of the above lakes. These 87 datasets of lake water
14 205 quality were obtained from a combination of sources including government, tribal agencies, university
15 206 researchers, citizen scientists, and non-profit agencies. Samples were taken during any season of the year
16 207 from the most recent decades, mostly from the late 1980's to 2012.

17 208 The largest challenge in building LAGOS-NE was the heterogeneity of the dataset formats,
18 209 variable conventions and units, and metadata, none of which were standardized. Many steps of data
19 210 integration required manual input from experts in diverse fields and close collaboration among specialists
20 211 in ecoinformatics, database design, freshwater ecology, and geography; all combined, the effort took six
21 212 years and involved ~15 individuals, spread across numerous institutions.

22 213 We designed the database using principles of open science so futures users could ask new
23 214 research questions by using the existing database or adding new data modules to the database. To ensure
24 215 users could do this, we documented the major steps of dataset integration and carefully integrated
25 216 metadata directly into the database itself, we emphasized data provenance, and we used a database
26 217 versioning system. In this data paper, we make the following research products available: (1) data tables
27 218 with the data that make up LAGOS-NE and an R package for accessing the data and integrating the
28 219 tables; (2) for each of the 87 water quality datasets, we provide the EML (ecological metadata language)
29 220 metadata files that we authored after receiving the data, the data files that we processed to import into
30 221 LAGOS-NE, and the R-script that we wrote to process the data; and (3) GIS coverages of the underlying
31 222 freshwater geographic features (lakes, streams and wetlands) that are linked to the data tables for GIS
32 223 processing by researchers.
33 224

34 225 **2. Study site: Midwest and Northeast U.S. lakes**

35 226 We selected an area of the U.S. known to have large numbers of lakes, well-developed lake water
36 227 quality sampling programs, and that spans diverse geographic conditions and thus gradients of ecological
37 228 context (Table 1). Our study area of 17 U.S. states includes 51,101 lakes ≥ 4 ha (Figure 1). These states
38 229 are in the north temperate climatic zone, which experience cold winters and warm, humid summers. The
39 230 study area includes part of the Interior Plains, Laurentian Uplands, Appalachian Highlands, and Atlantic
40 231 Plain geological provinces, and thus encapsulates a range of geological ages, glacial histories, and
41 232 topography. Land use/cover is highly variable, ranging from regions of intense agriculture in the corn belt
42 233 that spans portions of Minnesota, Wisconsin, Iowa, Missouri, Indiana, and Ohio, to predominantly
43 234 forested or urban regions of the northeastern U.S., including the states of Maine, New Hampshire, New
44 235 Jersey, and parts of New York, and primarily forested regions of northern Minnesota, Wisconsin, and
45 236 Michigan.

46 237 Although the majority of the data that we provide are for lakes ≥ 4 ha (see below for reasons for
47 238 using this threshold), we do include some data on lakes ≥ 1 ha and < 4 ha if data were available. Although
48 239 there may be water quality data for some lakes in this smaller size range, ecological context variables are
49 240 not available for these lakes.
50 241
51 242
52 243
53 244
54 245

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

State	Area (km ²)	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	2,819	11.3	1005	68.9	11.9	15.0	1.7
Indiana	93,717	1,874	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	2,645	5.1	1149	3.7	3.5	66.9	12.1
Massachusetts	21,013	1,698	8.9	1235	5.8	25.2	50.1	12.2
Michigan	150,489	6,511	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	1,858	12.7	1100	50.7	7.0	36.6	2.1
New Hampshire	23,980	1,109	6.5	1209	3.8	7.9	74.5	6.4
New Jersey	19,599	1,143	11.8	1188	13.8	31.1	27.9	21.4
New York	126,070	4,461	7.6	1094	21.9	9.3	54.1	7.2
Ohio	106,917	1,279	10.6	1003	50.0	14.7	30.9	1.0
Pennsylvania	117,293	1,755	9.3	1109	22.7	12.3	59.5	1.6
Rhode Island	2,809	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	6,009	6.6	831	36.7	7.5	35.5	13.7

This table includes 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 normals (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 one state.

3. 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 three main categories of variables: (1) variables that describe the physical characteristics and location of lakes themselves; (2) variables that describe in-situ water quality; and (3) 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 [18,19,20,12]. 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, or 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 three 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 (Figure 2). The module structure also facilitates data reuse and extension by accommodating future data modules related to any other lake or ecological-context feature.

1
2
3
4 280
5 281 **Figure 2. LAGOS-NE data modules and version numbers.** The data modules and versions that are included in
6 282 LAGOS-NE and are available with this paper include: LAGOS-NE_{GEO} v.1.05, LAGOS-NE_{LOCUS} v.1.01 (note, that in
7 283 Soranno et al. [17], this module was called LAGOS-lakes), and LAGOS-NE_{LIMNO} v.1.087.1. We include descriptions
8 284 of the type of data that are included in each module; with the major categories of variables the same as those
9 285 describing the data tables in Additional File 1. The black connectors among the modules show that the modules are
10 286 connected to each other through common unique identifiers through the LAGOS-NE_{LOCUS} module (through the
11 287 unique lake ID). P is phosphorus, N is nitrogen, C is carbon, S is sulfur, atm is atmospheric, NHD is the National
12 288 Hydrography Dataset, IWS is the interlake watershed, WBD is the Watershed Boundary Dataset, EDU is Ecological
13 289 Drainage Unit. Figure is modified from Figure 1 in Soranno et al. [17].
14 290
15 291

16 292
17 293 The design of LAGOS-NE and the workflow for its construction have been described previously
18 294 in detail [17]. In particular, the database design is based on the Consortium of Universities for the
19 295 Advancement of Hydrologic Science, Inc. (CUAHSI) Community Observations Data Model (ODM)
20 296 (CUAHSI ODM) as described in [17]. Here, we provide a brief overview. One important guiding
21 297 principle in creating LAGOS-NE was to ensure data provenance, i.e., that we could trace the original
22 298 source data through to the final LAGOS-NE database. Because each data module had different types of
23 299 source data, we developed different procedures for data provenance for each module, described in
24 300 Soranno et al. [17] and in this paper. The database model is based on ODM because it is a flexible data
25 301 model (i.e., allows the incorporation of wide range of types of data) that allows for the incorporation of
26 302 controlled vocabulary and, importantly, allows for extensive documentation through a relational database
27 303 structure of linked tables containing metadata [17]. The database was created and is maintained in
28 304 PostgreSQL v9.1. However, for researchers to use the database for analysis and modeling, it is necessary
29 305 to export the data into tables that can be processed by statistical packages or computer code. Therefore,
30 306 we exported the data into a series of tables (of similar data) that are needed to conduct research on either
31 307 the census population of lakes, the lakes for which there are water quality data, or some combination.
32 308 These are the data files that have been used to conduct research on LAGOS-NE to date, and that we make
33 309 available in this data paper (see Additional File 1 for a list of the tables and associated data that we are
34 310 making available). Further, we also make our GIS datasets available to facilitate geospatial analyses of
35 311 lakes, streams, and wetlands used to create some of the major components of LAGOS-NE.
36 312
37 313

40 314 **4. Description of LAGOS-NE_{LOCUS} v1.01 data module**

41 315 The LAGOS-NE_{LOCUS} module includes data on the physical location, some features and unique
42 316 identifiers for all lakes in the study area ≥ 1 ha, which means this data file has information on 141,378
43 317 lakes. Note, that because we detected errors in the digitization of lakes between 1 and 4 ha, we have
44 318 chosen to define our census population of lakes as only those ≥ 4 ha, but we still make data available for
45 319 lakes smaller than 4 ha when available in this and the LAGOS-NE_{LIMNO} data module. However, we
46 320 recommend caution in analyses, interpretation, and inference for lakes < 4 ha in this database that depend
47 321 on NHD's spatial representation and detection of water bodies. The data in this module include: lake
48 322 unique identifiers, perimeter, area, latitude and longitude (which is typically the centroid of the lake or a
49 323 central point that is within the lake boundary), GNIS name, and the zone IDs that the lake is located
50 324 within (e.g., state, county, or hydrologic units). The GIS datasets that we also make available provide the
51 325 lake polygon features associated with this module, as well as coverages for: lake watersheds, streams,
52 326 wetlands, spatial classifications, and glaciation history.

53 327 Definition of lakes: We defined lakes previously in Soranno et al. [17] as follows. A 'lake' in
54 328 LAGOS-NE is a perennial body of relatively still water. We include lakes and reservoirs that range from
55 329 being completely natural to highly modified: lake basins can be entirely natural, modified natural (i.e., a
56 330 water control structure on a natural lake), or a fully impounded stream or river (i.e., a reservoir). We
57 331 explicitly exclude: sewage treatment ponds, aquaculture ponds, and detention ponds that are known to
58 332
59 333
60
61
62
63
64
65

1
2
3
4 331 contain basins that are entirely artificial and were built for high-intensity human use. In addition, due to
5 332 their unusual nature and size, we do not include the five Great Lakes in our database. This definition of
6 333 'lake' for LAGOS-NE has been developed only for the purpose of this database and its applications (e.g.,
7 334 to answer questions about lake water quality). The intent of LAGOS-NE is not to document and measure
8 335 the total number of water bodies in our study area, although we are able to perform this calculation for
9 336 lakes ≥ 4 ha, with an acceptable level of uncertainty (see below).

10 337 Definition of lake watersheds: We calculated lake watersheds as 'inter-lake watersheds' (IWSs)
11 338 defined as the area of land draining directly into the lake as well as the area that drains into upstream-
12 339 connected streams and lakes < 10 ha (Figure 3). We defined lake watersheds this way to define the
13 340 drainage basin of lakes that includes connected streams and their drainage basins. However, because
14 341 research has shown that large upstream lakes can trap nutrients flowing into them, these large lakes can
15 342 block nutrient transport of nutrients that originate upstream from them to downstream lakes in a
16 343 connected lake chain (e.g., [22]). Therefore, to calculate a drainage basin for a lake with large upstream
17 344 connected lakes, we did not include the drainage basins of upstream lakes > 10 ha. See Soranno et al. [17]
18 345 for full details on how lake IWSs were calculated and the section on LAGOS-NE_{GEO} for further details.

19 346 Lakes near and beyond the state borders: For some of our analyses, we delineated boundaries in
20 347 other ways than political boundaries that were more ecologically relevant, which resulted in the inclusion
21 348 of some lakes outside of the exact 17 state border. This fact allowed us to include more in situ data
22 349 collected by state and citizen sampling programs which do not always follow strict state borders and may
23 350 include lakes that are outside of state lines. Although most of these border lakes have hydrological (i.e.,
24 351 lake connectivity measures) and topographic (i.e., lake watershed delineations) calculations or water
25 352 quality data, some measures of ecological context may be missing. For example, for lakes in Canada, we
26 353 were not able to estimate any data that relied on national datasets that stopped at the Canadian border; one
27 354 exception is the NHD, which extends into Canada to retain hydrologic boundaries.

31 355 **Data sources of the LAGOS-NE_{LOCUS} module**

32 356 Detailed information on data sources are found in 'Additional File 5' in Soranno et al. [17]. Briefly, the
33 357 data source for lakes and streams in the 17 state area was the NHD [21]. The hydrologic boundaries (i.e.,
34 358 for three of the spatial classifications, HUC12, HUC8, HUC4) came from the Watershed Boundary
35 359 Dataset (WBD; [23]). In addition, we used the digital raster dataset of elevation for watershed delineation
36 360 from the National Elevation Dataset [24]. All download dates for these data sources are provided in
37 361 'Additional File 5' in the above citation.

38 362 **Data-integration methods of the LAGOS-NE_{LOCUS} module**

39 363 All methods to create this module are described in Soranno et al. [17]. The most challenging and time-
40 364 consuming part of building this module was connecting the sampling locations from the lake water
41 365 quality datasets (which each contained different types of unique identifiers, and sometimes only lake
42 366 names) to a georeferenced location in the NHD. When data providers included the lake latitude and
43 367 longitude, we were able to mostly automate the procedure. Nevertheless, even when coordinates were
44 368 available, there were many cases where the latitude and longitude did not intersect the NHD lake polygon
45 369 boundary, requiring manual interpretation.

46 370 **Figure 3. Examples lake watersheds (IWS) in LAGOS-NE.** The watersheds are coded by hydrologic class to
47 371 which its lake belongs. Data are from the LAGOS-NE_{GEO} v.1.01 data module and the GIS data coverages.

48 372 **Quality Control of the LAGOS-NE_{LOCUS} module**

49 373 The full description of error analysis for this module is described in Soranno et al. [17]. However,
50 374 here we briefly describe our efforts to determine the minimum area of a lake that we could confidently
51 375
52 376
53 377
54 378
55 379
56 380
57 381

represent using the NHD (further details located in Additional File 9 in Soranno et al. [17]). Although the NHD is a national dataset, 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, 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 four states (WI, MI, IA, ME) in which to evaluate error rates of water body identification for lakes ≥ 1 ha and seven states (WI, MI, IA, ME, MO, NH, OH) in which to evaluate error rates for lakes ≥ 4 ha. We randomly selected three 100-km² rectangles from each state then compared the number of lakes occurring in the NHD GIS coverage to 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 ≥ 1 ha four-state test and 13% for the ≥ 4 ha seven-state test. Because an average of 87% of lakes ≥ 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_{LOCUS} module

Figure 1 shows the census population of all lakes ≥ 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 (9,808 lakes with maximum depth values, and 4,090 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 (Figure 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).

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

State	Lakes ≥ 4 ha (#)	Isolated Lakes (#)	Headwater lakes (#)	Drainage lakes (#)	Drainage lakes with upstream lakes (#)
Connecticut	770	40	119	424	187
Illinois	2,831	1,417	279	952	183
Indiana	1,883	760	244	697	182
Iowa	915	339	87	402	87
Maine	2,661	94	619	1,211	737
Massachusetts	1,716	210	269	751	486
Michigan	6,531	2,649	1,087	1,672	1,123
Minnesota	14,031	6,609	1,894	2,673	2,855
Missouri	1,865	435	179	1,113	138
New Hampshire	1,118	70	224	581	243
New Jersey	1,148	219	129	521	279
New York	4,477	629	1,210	1,915	723
Ohio	1,282	543	105	520	114
Pennsylvania	1,757	316	397	840	204
Rhode Island	266	35	40	115	76
Vermont	531	14	74	364	79
Wisconsin	6,026	2,982	823	1,236	985
Total	49,808	17,361	7,779	15,987	8,681

1
2
3
4 417 The number (#) of lakes ≥ 4 ha in each of the lake hydrologic classes by state, as well as the total numbers of lakes
5 418 by hydrologic class calculated for the study extent. Note, in this table, lakes are counted for each state in which they
6 419 occur (i.e., lakes that straddle two states are counted in both states).
7
8

9 421 **5. Description of LAGOS-NE_{LIMNO} v1.087.1 data module**

10 422
11 423 The LAGOS-NE_{LIMNO} module includes in situ measurements of lake water quality. We included
12 424 variables that are most commonly measured by state agencies and researchers for studying eutrophication
13 425 (Figure 2, **Water quality data and metadata**, including chlorophyll a, Secchi depth, water color, DOC,
14 426 total and dissolved phosphorus (P), nitrogen (N), and carbon (C)). For each water quality data value, we
15 427 also include metadata as additional columns in the exported data table (Figure 2, **Metadata** including
16 428 analytical method, data qualifier from the program, detection limit (when available), and the LAGOS-NE
17 429 sensor code) including: the analytical methods, qualifiers with data flags from the original program (*qual*,
18 430 which is not standardized for LAGOS-NE), detection limits (if available), and standardized sensor codes
19 431 from our quality control procedures (*sensorcode*, standardized for LAGOS-NE). Finally, we include
20 432 documentation about each source program that is linked to each data value.
21 433

22 433 23 434 Data sources of the LAGOS-NE_{LIMNO} module

24 435 We acquired individual water quality datasets for LAGOS-NE_{LIMNO} by contacting individuals at
25 436 each of the 17 state and 5 tribal agencies. These contacts helped us to identify the state-agency collected
26 437 dataset required by the Clean Water Act and which is most likely to be in the public domain. In this way,
27 438 we were able to acquire at least one (and typically more) dataset from each of the 17 states. Because state
28 439 and tribal agencies vary in sampling approach and intensity (see below for details), we sought to
30 440 supplement these datasets with other known sources of water quality data, including university
31 441 researchers, federal agencies, and non-profit groups to integrate into the LAGOS-NE_{LIMNO} module. The
32 442 full list of data sources acquired is in Soranno et al. [17] in ‘Additional File 17’; however, we
33 443 incorporated a subset of these datasets in LAGOS-NE_{LIMNO} v1.087.1 (the data file
34 444 *LAGOSNE_source_program_10871.csv* contains the list of sources for this version of LAGOS-NE).
35 445

36 445 37 446 Data-integration methods of the LAGOS-NE_{LIMNO} module

38 447 All methods to create this module are described in Soranno et al. [17]. Briefly, for each dataset acquired,
39 448 we authored LAGOS-NE metadata in EML to aid in data provenance (included in this paper). We also
40 449 incorporated key metadata features (e.g., methods used, sensor codes (if applicable)), and sampling
41 450 program information) into the database so that future users could easily identify these important
42 451 attributes. Because each dataset was unique in structure, file format, and naming conventions, we
43 452 manually processed each dataset and its metadata so that they could be translated into the standard
44 453 LAGOS-NE vocabulary and data model. Although labor-intensive, we created customized R scripts to
45 454 process and load each dataset separately (included in this data paper).
46 455

47 455 48 456 Quality control of the LAGOS-NE_{LIMNO} module

49 457 The full description of our quality assurance/quality control (QAQC) procedures for this module are
50 458 described in Additional File 2. Here, we provide a brief overview of our approach. Our goal for this effort
51 459 was to identify egregiously high values and values that might be too low, both defined below. Note that
52 460 our quality control procedures were not designed to identify statistical outliers, which individual users are
53 461 expected to perform themselves because such analyses depend on the subsequent statistical analysis of
54 462 each user. There were three major phases in the QAQC procedure for LAGOS-NE_{LIMNO}. Phases I and II
55 463 were designed to identify the egregious values that we defined as those that: (1) did not make ecological
56 464 sense, (2) were far beyond what has been detected in previous studies, (3) were not technically feasible
57 465 (e.g., SRP > TP), or (4) were a result of a data or file corruption or error in the data loading stage. For
58 466 these egregious values, we explored the issues that might be underlying the values and removed them
59 467 from the LAGOS-NE_{LIMNO} data export provided in this data paper because we had sufficient evidence that
60
61
62
63
64
65

1
2
3
4 468 they were not scientifically valid data values. We were very conservative in these assessments to avoid
5 469 removing data values that were high, yet still valid. Phase III was designed to identify and flag values that
6 470 seemed to be lower than analytically possible (i.e., below detection limits) when there was sufficient
7 471 metadata; however, note that these data are still provided in this data paper because it is not appropriate to
8 472 remove data that are below detection when those data could be valid.

9 473 For all versions of LAGOS-NE_{LIMNO}, Phase I and II are conducted on the entire cumulative
10 474 dataset to leverage as large of a sample size as possible to detect problem values. In other words, because
11 475 many of the QAQC analyses outlined here make use of all information from an individual lake or
12 476 variable, incorporating new data may result in a better assessment of the data than when there are less
13 477 data. Thus, for each new version of LAGOS-NE_{LIMNO}, new decisions are made about egregious values. In
14 478 this data paper, we describe the procedures for assessing all major versions of LAGOS-NE_{LIMNO}, but we
15 479 present the results only for this version of LAGOS-NE_{LIMNO} (v1.087.1).

16 480 Because there are few accepted practices for conducting such quality control on a large,
17 481 integrated database, we created our own procedures for Phase I and II by creating tests to identify
18 482 egregious values that leverage a large, integrated database with multiple measures of water quality and
19 483 well-established expected relationships among variables. The database that we used to identify egregious
20 484 values was based on data in the full LAGOS-NE_{LIMNO} database for samples taken from all lake depths
21 485 provided by the source datasets (note, our data exports in this data paper are only for epilimnetic or
22 486 surface samples). While the quality control procedures that we implemented here were designed to help
23 487 resolve the large and egregious errors in a combined dataset such as this, there are likely additional
24 488 extreme values in the database due to the size and heterogeneity of the data. Users may want to check for
25 489 additional issues in the data values specific to their intended analyses.

26 490 27 491 **Data in the LAGOS-NE_{LIMNO} module**

28 492 All data in LAGOS-NE_{LIMNO} v1.087.1 are from samples that we identified as being collected
29 493 from either the lake surface or the epilimnion (the well-mixed surface layer of a thermally-stratified lake
30 494 during the period of stratification). Because we did not have lake temperature data to quantify the exact
31 495 epilimnion depth in all lakes, we used information from the source datasets to either determine epilimnion
32 496 depth, or to select data from only the top water layers. Although we received data from different depths in
33 497 lakes, the majority of the samples were from the surface or epilimnion. The database includes samples
34 498 from any season of the year. However, most of the published analyses to date have focused on the
35 499 summer stratified period.

36 500 Lakes are not sampled the same way by all individuals, groups, or agencies; there are differences
37 501 in the variables measured, the frequency and timing of sampling, and the proportion of lakes sampled. For
38 502 example, for total phosphorus, the four states with the largest number of unique lakes with at least one
39 503 value for total phosphorus per state include: Wisconsin (1,920 lakes), Minnesota (1,588), New York
40 504 (1,289), and Michigan (1,109) (Table 3). However, the states with the highest proportion of their lakes
41 505 with total phosphorus samples are the smaller states with fewer numbers of lakes, such as New
42 506 Hampshire (64%), Vermont (58%), and Rhode Island (42%). Notably, there are some states with
43 507 intermediate numbers of lakes that still have quite large percentages of their lakes with total phosphorus
44 508 values, including Maine (35% of 2,645 lakes), Wisconsin (32% of 6,009 lakes), and New York (29% of
45 509 the 4,461 lakes).

46 510 The most commonly measured variable in LAGOS-NE_{LIMNO} is water clarity measured as Secchi
47 511 depth (a relatively easy and cost-effective measure of water quality), with 897,724 measurements taken
48 512 from 12,034 unique lakes in the 17 states from mostly the mid 1980s to 2011 (Table 3). The second and
49 513 third most sampled measures of water quality are chlorophyll *a* and total phosphorus, respectively.
50 514 Although it appears that total nitrogen is sampled far less frequently than total phosphorus, some labs
51 515 measure total nitrogen directly and report that single value, whereas other labs measure the constituents
52 516 that make up total nitrogen (total Kjeldahl nitrogen and nitrate+nitrite), and sum them together to
53 517 calculate total nitrogen. All of our analyses conducted on total nitrogen have used such calculated and
54 518 measured values of nitrogen together, which increase the sample sizes for total nitrogen markedly.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

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

14
15
16
17
18
19
20 522
21 523
22
23
24

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

State	Number of lakes (≥4 ha)	Variable	Total phosphorus	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
		sample years:	1972-2010	1937-2010	1937-2013	1984-2007	n/a	1984-2007	1973-2010	1999-2009	1976-2010
Illinois	2819	# of samples:	2816	2317	1438	20	0	20	43	1526	2351
		# 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
Indiana	1874	# of samples:	1232	1303	909	57	0	57	57	1183	1237
		# of sampled lakes:	341	340	320	51	0	51	51	322	341
		sample years:	1988-2010	1986-2010	1990-2009	2007	n/a	2007	2007	1988-2009	1988-2009
Iowa	903	# of samples:	2873	2836	2711	18	0	18	2244	6	2229
		# of sampled lakes:	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
Maine	2645	# of samples:	17314	83472	12480	1927	1676	3321	1260	8	1577
		# of sampled lakes:	933	1047	793	601	466	848	461	3	347
		sample years:	1971-2011	1952-2011	1974-2011	1983-2011	1972-2011	1984-2011	1995-2011	1978-1993	1978-2011
Massachusetts	1698	# of samples:	570	760	326	277	228	300	69	69	351
		# 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
Michigan	6511	# of samples:	10143	95283	12243	1811	69	987	749	2651	4850
		# of sampled lakes:	1109	1233	862	836	69	353	200	713	948
		sample years:	1965-2013	1925-2013	1959-2013	1973-2010	2002-2003	1984-2013	1959-2011	1980-2010	1973-2012
Minnesota	13984	# of samples:	10974	497646	81925	406	6683	3382	7717	43054	7725
		# 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
Missouri	1858	# of samples:	11619	11794	11578	27	0	27	11340	0	27
		# of sampled lakes:	208	207	201	23	0	23	207	0	23
		sample years:	1978-2013	1978-2013	1978-2013	2007	n/a	2007	1978-2013	n/a	2007
New Hampshire	1109	# of samples:	9289	2958	154	237	3044	390	22	1209	2445
		# 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	1975-1994	1975-2013
New Jersey	1143	# of samples:	421	461	446	27	0	44	10	443	472
		# of sampled lakes:	175	174	157	25	0	36	8	157	175
		sample years:	1984-2009	1984-2009	2005-2009	1984-2007	n/a	1984-2007	2007	2005-2009	1984-2009
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
		sample years:	1975-2012	1975-2012	1975-2012	1981-2012	1984-2011	1982-2011	1990-2012	1981-2010	1975-2012
Ohio	1279	# of samples:	377	1868	1912	20	0	220	1873	0	447
		# 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
Pennsylvania	1755	# of samples:	1170	924	971	163	0	160	638	16	290
		# of sampled lakes:	263	260	160	124	0	124	167	2	147
		sample years:	1980-2011	1984-2011	1980-2011	1984-2008	n/a	1984-2007	1997-2011	1985-2010	1980-2010

61
62
63
64
65

14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Rhode Island	253	# of samples:	3325	18211	12195	51	6	65	2582	0	2100
		# 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
Vermont	528	# 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
Wisconsin	6009	# 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
TOTAL	49592	# 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 one 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_{LIMNO} 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 nitrogen, total dissolved phosphorus, total organic carbon, and total organic nitrogen. n/a is not applicable.

14
15
16
17
18
19
20 535
21 536
22 537

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

Program Type	Number of datasets	Number of lakes (≥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	# of values:	419	527	324	229	173	215	335	6	30
			# of unique lakes:	17	17	17	13	15	14	16	1	9
Federal Agency/ University	2	2	# of values:	-	799	-	-	-	-	-	-	-
			# of unique lakes:	-	2	-	-	-	-	-	-	-
LTER	3	9	# of values:	2,346	3,529	2,567	-	-	1,872	1,612	507	2,396
			# of unique lakes:	9	9	5	-	-	9	9	4	9
National Survey Program	5	2,244	# of values:	2,320	2,595	243	3,689	703	4,714	431	-	4,204
			# of unique lakes:	1,863	1,891	171	13	142	2,235	398	-	1,997
Non-Profit Agency	4	44	# of values:	1,326	4,798	2,678	-	-	-	214	9	908
			# of unique lakes:	44	41	28	-	-	-	39	1	44
State Agency	33	4,264	# of values:	34,348	42,888	29,993	16,240	5,010	14,528	5,359	7,220	25,684
			# of unique lakes:	3,914	3,186	2,309	2,092	776	1,191	634	1,991	3,216
State Agency/ Citizen Monitoring	11	7,039	# of values:	79,390	645,650	124,766	18,010	8,630	3,195	18,610	52,995	27,826
			# of unique lakes:	3,955	6,629	4,341	1,111	1,508	786	772	3,476	2,782
State Agency/Univ/ Citizen Monitoring	4	1,835	# of values:	31,809	194,177	37,993	439	1,171	1,519	10,844	-	2,112
			# of unique lakes:	1,439	1,812	1,253	302	393	574	712	-	99
Tribal Agency	5	46	# of values:	911	145	905	3	-	357	411	277	463
			# of unique lakes:	33	3	32	3	-	11	18	5	17
University	17	535	# of values:	2,273	4,412	3,939	172	69	723	2,275	-	2,397
			# of unique lakes:	326	500	415	151	69	318	396	-	171

53 538
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 539
5 540 **Figure 4. Percentage of lakes by lake area with water quality data.** Percentage of census lakes in each lake size
6 541 bin (top panel) compared to the percentage of census lakes for which there are limnological data for Secchi (second
7 542 panel), chlorophyll *a* (third panel), and total phosphorus (TP; bottom panel).
8 543

9 544 Using the three most sampled variables in the dataset (Secchi depth, chlorophyll concentration
10 545 and total phosphorus), we found that larger lakes were more likely to be sampled for water quality than
11 546 smaller lakes (Figure 4). This result was expected given the economic and recreational interest in larger
12 547 lakes, including easier public access. Previous research has already documented this basic pattern in 6 of
13 548 the states included in LAGOS-NE [30]. Across all states, almost 80% of lakes > 400 ha have water
14 549 quality data.

15 550 Lakes are also unevenly sampled through time, depending on the variable (Figure 5). Some
16 551 programs' focus is on long-term monitoring, whereas others are short-term initiatives. Typically, long-
17 552 term monitoring programs are localized to a few lakes, although there are exceptions (e.g., monitoring for
18 553 acid rain in the northeastern US in the 1980s-present has resulted in good temporal and spatial coverage
19 554 for some variables through time and space [31].
20 555

21 556 **Figure 5. The number of years of water quality data by lake.** The number of years for which at least one sample
22 557 is taken during the summer stratified season (15 June to 15 September) for: Secchi depth in meters, total phosphorus
23 558 in µg/L, total nitrogen in µg/L (includes both measured and calculated values), and chlorophyll *a* in µg/L.
24 559
25 560

26 561 **6. Description of LAGOS-NE_{GEO} v1.05 data module**

27 562 The LAGOS-NE_{GEO} module includes information on the ecological context of the census lakes,
28 563 their watersheds, and their regions. The information provided in the data tables for this module is
29 564 organized into three main themes in which data are exported into individual tables: CHAG - climate,
30 565 hydrology, atmospheric deposition of nitrogen and sulfur, and surficial geology; LULC - land use/cover,
31 566 canopy cover, terrain metrics, and dam density; and CONN - lake, stream, and wetland abundance and
32 567 connectivity measures (Figure 2). We also provide the GIS coverages that include some of the underlying
33 568 data for this module, including: lake polygons and their hydrologic classifications defined in [17];
34 569 wetland polygons and their classification; streams as a line coverage and their classification by stream
35 570 order; the zones used for this study (state and county; hydrologic units [at the 4, 8 and 12 scales; [32]]);
36 571 and, lake watersheds (IWS). We also include boundaries of U.S. states and Canadian provinces for
37 572 mapping.
38 573

39 574 **Data sources of the LAGOS-NE_{GEO} module**

40 575 Detailed information on data sources are found in 'Additional File 5' in Soranno et al. [17].
41 576 Almost all data sources for this module are from national-scale datasets and thus use standardized
42 577 methods throughout the study extent.
43 578

44 579 **Data-integration methods of the LAGOS-NE_{GEO} module**

45 580 All methods to create this module are described in 'Additional files 5, 7, 8, 13, and 14' in
46 581 Soranno et al. [17]. Briefly, we calculated the metrics for this module that describe the ecological context
47 582 surrounding lakes by developing project-specific GIS tools in the ArcGIS environment, which are
48 583 referred to as the LAGOS GIS Toolbox (and made available here: [33]). The toolbox outputs multiple
49 584 individual data tables of calculated values organized by the above three data themes that are then
50 585 imported into LAGOS-NE_{GEO} for different spatial classifications, including values calculated at the level
51 586 of the individual lake, 100-m and 500-m buffers around each lake, the lake IWS, states and counties,
52 587 hydrologic units, and ecological drainage units (an ecoregion spatial classification). The unique identifiers
53 588 for this data module are the zone ID's for each spatial classification for which we calculate these metrics.
54 589 In other words, we calculate land use around a lake in each of the zones of the many spatial classifications
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 590 in LAGOS-NE. However, the data are exported into individual tables by spatial classification. Therefore,
5 591 there are different numbers of rows in each table; for example, there are 51,101 rows for the land use
6 592 metrics calculated for the 100-m lake buffer because there are 51,101 lakes that have a 100-m buffer area,
7 593 but only 17 rows for the land use metrics calculated for the state spatial classification.
8
9 594

10 **Quality control of the LAGOS-NE_{GEO} module**

11 596 The full description of error analysis for this module is described in ‘Additional file 14’ in Soranno et al.
12 597 [17]. The quality control procedures for this module included procedures to identify possible errors or
13 598 improbable values as a result of the extensive automated GIS data processing that creates the LAGOS-
14 599 NE_{GEO} data tables and to correct those problems. We assumed that the original data layers had already
15 600 gone through extensive quality control by the originators of the datasets. We defined errors and
16 601 improbable values to be: (1) values that did not make ecological sense; (2) values that were well beyond
17 602 what has been observed in previous studies; (3) values that are not technically feasible; or, (4) null values
18 603 that indicate an absence of data, when in fact data exist based on the input data coverages. Note, it was
19 604 not our intention to remove statistical outliers that may or may not be real/true values. Rather, we
20 605 conducted procedures on each exported table that included: verifying column headers and units, mapping
21 606 the exported data to evaluate mapping extent and boundary issues using visual inspection, mapping the
22 607 data distributions of each value, identifying values that were missing or zero, plotting distributions of the
23 608 data, ensuring that proportions summed to 100 where relevant, and inspecting univariate plots of metrics
24 609 that are known to be related (e.g., % urban land use versus % impervious surface).
25
26 610

27 **Data in the LAGOS-NE_{GEO} module**

28 611 This module contains the largest amount of data of any of the modules. For example, Figure 6 shows the
29 612 wide range of ecological context for the LAGOS-NE study area calculated for three different spatial
30 613 classifications. For those variables that are measured coarsely (e.g., baseflow, runoff, atmospheric
31 614 deposition, geology), we calculated variables for only the broader spatial classifications. For example, we
32 615 did not calculate baseflow for spatial classifications finer than HUC12 because the underlying data for
33 616 baseflow is estimated on a zone generally coarser than the area of a lake watershed.
34
35 617
36 618
37 619
38 620

39 **Figure 6. Example ecological context variables by spatial classification in LAGOS-NE.** The top four panels are
40 622 zoomed in to selected regions of Minnesota and Wisconsin so that the zone boundaries can be seen. The upper left
41 623 panel shows stream density in each lake IWS, and the upper right panel shows the percent of connected wetlands in
42 624 each lake IWS. The middle left panel shows the 2011 percent urban land use/cover in each hydrologic unit code 12
43 625 (HUC12), and the middle right panel shows the 2011 percent agricultural land use/cover in each hydrologic unit
44 626 code 12 (HUC12). The lower left panel shows the 2010 nitrogen deposition in each HUC8, and the lower right panel
45 627 shows the average percent of streamflow that is baseflow in each HUC8.
46 628
47 629
48 630
49 631
50 632
51 633
52
53
54
55
56
57
58
59
60
61
62
63
64
65

7. 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 its wide gradients of ecological context. The lake census dataset also makes it possible to quantify the types of biases present in the dataset 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 three 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 two 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 two 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: (1) a lack of lake depth information (lake depth drives many in-lake processes), (2) 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, (3) 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] has 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 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. (in prep) 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_{LIMNO} 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. (in prep) 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_{LIMNO}, including v1.087.1.

1
2
3
4 683 Lake water quality is affected by many ecological context features, such as lake physical
5 684 characteristics, land cover, land use, and climate. The relationship between these features and the water-
6 685 quality measurements is not always linear. In addition, the data tend to be noisy and often contain missing
7 686 values, which makes it challenging to fit effective statistical models. To overcome these challenges, Yuan
8 687 et al. [37] developed a novel algorithm for learning non-linear features to predict lake water quality. The
9 688 algorithm also enables the missing values to be imputed in a way that preserves the relationship between
10 689 the predictors and response variables. Furthermore, because many of the lake water-quality variables are
11 690 strongly correlated with each other, their models are expected to be similar. This similarity information
12 691 can thus be exploited to build better models especially for the lake water-quality variables that have very
13 692 few observations because they are not sampled frequently. Yuan et al. (in prep) are developing a machine
14 693 learning approach known as multi-task learning that can simultaneously build regression models of
15 694 multiple lake water-quality variables for a large number of lakes, taking into account both the correlation
16 695 between the variables and the spatial autocorrelation among the lakes. Because we expect many
17 696 ecological datasets across broad geographic scales to have similar data gaps and challenges as LAGOS-
18 697 NE, we think these methods will be extremely valuable for other researchers studying different
19 698 macroscale questions.
20 699

23 700 *Understanding spatial variation in lake nutrients and eutrophication at sub-continental scales:*

24 701 LAGOS-NE allows investigation of spatial variation in lake nutrients and eutrophication at
25 702 macroscales. For example, Lapierre et al. (in prep) identify general spatial principles that constrain
26 703 relationships between ecosystem variables with different spatial structures. In other cases, specific
27 704 questions regarding spatial patterns have focused on identifying important landscape controls on nutrients
28 705 and their ratios [38], potential stress induced on phytoplankton communities by high nitrogen levels
29 706 (Filstrup et al. in prep), and spatial autocorrelation in lake-specific relationships between chlorophyll and
30 707 nutrients and carbon [39]. In addition, LAGOS-NE contains a wealth of information on a variety of lake
31 708 ecosystem types. Shallow lakes, in particular, are very abundant across the study area and represent
32 709 systems that can exhibit hysteresis in response to lake eutrophication. Cheruvelil and Wagner (in prep) are
33 710 investigating the spatial distribution and temporal dynamics of water clarity in shallow lakes of the
34 711 LAGOS-NE study area.

35 712 An important area of research, and one that was a motivating factor for the creation of LAGOS-
36 713 NE, is understanding the importance of cross-scale interactions (CSIs) — where ecological processes
37 714 operating at one spatial or temporal scale interact with processes operating at another scale — in lake
38 715 ecosystems. Because of their importance ecologically and the challenge of quantifying them over large
39 716 spatial extents, Wagner et al. [40] evaluated the statistical power of large multi-thematic, multi-scaled
40 717 datasets, such as LAGOS-NE, to detect CSIs. This work not only helped inform the design of large-scale
41 718 studies aimed at detecting CSIs, but also focused attention on the importance of considering CSI effect
42 719 sizes and their ecological relevance.. To extend this work, Fergus et al. (in prep) are investigating the
43 720 importance of both within- and cross-scale interactions in landscape models predicting lake nutrients, and
44 721 the role that connectivity among freshwaters plays in these interactions. Understanding and predicting
45 722 nutrients in lakes at macroscales is important to inform estimates of lake contributions to continental and
46 723 global nutrient cycles. To date, much of this work has been performed on a nutrient-by-nutrient basis,
47 724 despite knowing that cycles of nitrogen and phosphorus and other key elements are best understood by
48 725 considering multiple elements in tandem, *e.g.*, in a stoichiometric framework [41] or through analysis of
49 726 coupled biogeochemical cycles (*e.g.*, [42, 43, 44]). Currently, efforts are underway to develop spatial joint
50 727 nutrient distribution models to evaluate how our understanding of landscape-scale drivers of lake
51 728 nutrients and predictive performance are improved by considering multiple nutrients simultaneously
52 729 (multivariate models) compared to traditional univariate approaches that ignore that nutrient cycles can be
53 730 tightly coupled in freshwaters (Wagner et al. in prep).
54 731
55 732
56 733

1
2
3
4 734 Understanding temporal and spatial variation in lake eutrophication at sub-continental scales:

5 735 In addition to the vast spatial data contained in LAGOS-NE, temporal data are available for many
6 736 water-quality variables, and some of the ecological context variables (e.g., land use/cover and
7 737 atmospheric deposition). This is important information within the context of understanding and predicting
8 738 how lake ecosystems have and will respond to global change, such as changes in climate and land use,
9 739 and management activities to reduce nutrient inputs to lakes. Because we do not expect responses to such
10 740 change and actions to be the same everywhere, these questions must be addressed across both space and
11 741 time. In particular, recent environmental changes and management efforts have been hypothesized to both
12 742 improve and degrade water quality in lakes. However, to date, there have been no studies to examine
13 743 these issues comprehensively across broad scales and to examine which drivers are most strongly related
14 744 to eutrophication status in lakes. LAGOS-NE is very well suited to answer these types of questions.

15 745 For example, nearly 3,000 lakes were examined for trends in nutrients and chlorophyll from 1990
16 746 to 2013 using LAGOS-NE [45]. Across all lakes, nitrogen has declined, and phosphorus and chlorophyll
17 747 have not changed. Nitrogen and stoichiometric changes in lakes were related to atmospheric deposition of
18 748 nitrogen, providing key insight into large-scale nutrient transport and policies such as the Clean Air Act.
19 749 Using only citizen-science data in a subset of the LAGOS-NE database, Lottig et al. [46] showed results
20 750 that suggested little evidence for major declines or improvements in water quality. In addition, Collins et
21 751 al. (in press) are examining the relationships between a wide range of climate metrics and water quality in
22 752 ~11,000 lakes in LAGOS-NE to determine, 1) which climate metrics are most related to water quality; 2)
23 753 whether physical, chemical and biological aspects of lakes respond to climate in the same way; and, 3)
24 754 how the climate-water-quality relationship varies across space and regions with different ecological
25 755 context. However, the temporal dynamics of lake ecosystem properties can sometimes be nonlinear and
26 756 exhibit variability across the landscape--largely because of climate and within-lake processes. Lottig et al.
27 757 (in prep) have developed models for understanding and predicting the often complex temporal patterns
28 758 observed in water clarity. These studies point to the importance of considering both space and time when
29 759 trying to understand broad-scale environmental issues in surface waters.

30 760
31
32
33
34
35 761 **8. Using LAGOS-NE for future research, management, and policy**

36 762 To facilitate potential future use of LAGOS-NE, we have thoroughly documented the database
37 763 and its methods [17]; and, here, we share LAGOS-NE data with the broader research community. In this
38 764 data paper, we include a wide range of research products, including: the water quality and ecological-
39 765 context data; the GIS coverages underlying much of the analyses on freshwaters; and, an R package that
40 766 facilitates use of LAGOS-NE [47]. This package includes functions to retrieve, store, and interact with the
41 767 LAGOS-NE database that works across many different operating systems. The package should increase
42 768 the ease with which users of the database are able to access the data and documentation while maintaining
43 769 a reproducible workflow.

44 770 Key motives for constructing this database included interest in examining lake nutrients and
45 771 productivity at multiple spatial and temporal scales, fostering broad-scale aquatic ecology and
46 772 macrosystems research in an open-science platform, and providing new understanding and resources for
47 773 management and policymakers. To this end, several team members have made presentations at scientific
48 774 meetings about the structure and use of LAGOS-NE and subsets of LAGOS-NE data have been shared
49 775 with other researchers and stakeholders and agency personnel in advance of this publication. These early
50 776 uses of LAGOS-NE data by other researchers outside of our team include an investigation of patterns and
51 777 causes of shifting distribution of a sentinel fish species (Rypel et al. in prep), developing models to
52 778 simulate lake temperatures (Winslow et al. in prep) and fish species distributions, and developing a
53 779 recruitment model for a popular game fish (Hansen et al. in prep). Results from the latter two efforts will
54 780 inform state-level fisheries management as well as aid in prioritization of lakes for habitat conservation
55 781 action across a tri-state region.

56 782 Much of the research that we and others are conducting with LAGOS-NE has implications for
57 783 ecosystem management or environmental decision-making. In addition, we have collaborated with

1
2
3
4 784 boundary organizations and decision-makers. For example, under development is a dashboard of the
5 785 ecosystem services provided by lakes for use by land managers (Keeler et al. in prep). In addition, we
6 786 have helped the state of Michigan determine lake-specific nutrient standards [48]. Our hope is that this
7 787 database and the associated support tools and documentation serve as a powerful resource and a
8 788 foundation for future research and decision-making by a broad community of scientists, policy makers,
9 789 and natural-resource managers. Indeed, our success and experience with database construction and
11 790 research has inspired us to expand the spatial extent for LAGOS-NE. We have begun to build LAGOS-
12 791 US, which will include similar data as LAGOS-NE but will be for the continental U.S.
13 792

14 793 **9. Challenges and recommendations for creating large, integrated, and** 15 **heterogeneous databases** 16 794 17 795

18 796 We found that the largest challenge when creating this database was integrating many small
19 797 heterogeneous datasets that had few common standards. Although creating such large, integrated datasets
20 798 using fully automated procedures may happen someday, it appears that we are nowhere near such
21 799 automation today. Until standards in metadata documentation and robust ontologies are created and
22 800 widely adopted when creating local or regional datasets, future efforts to integrate these into larger
23 801 databases will have to rely on close collaborations among domain experts and ecoinformatics
24 802 professionals, extensive manual interpretation of individual datasets, and funds sufficient to implement
25 803 these labor-intensive approaches [16]. Nevertheless, it is worth the time and money invested in database
27 804 integration if the resulting databases support new research, management, policy, public outreach, and
28 805 education at all levels. We anticipate that LAGOS-NE will serve as a foundation for new data modules
29 806 that can be used beyond the original intent of LAGOS-NE.
30 807

31 808 The economic value of water quality data in an integrated database

32 809 This extensive effort was supported by a U.S. National Science Foundation grant that totaled \$2.4
33 810 million, along with resources from other projects. Our team ranged in size from 14-20 individuals across
34 811 the six years of the project, with many members compiling and integrating data, authoring metadata,
35 812 creating new data products, and implementing quality control procedures, resulting in a tremendous
36 813 number of person-hours. However, when one considers the cost of the data collection for the water quality
37 814 data in the first place, the expense of this post-processing integration work is not as large as it sounds.
38 815 Sprague et al. [16] suggest that a single sample (estimated for collecting nutrient or chemistry data from
39 816 streams) ranged in cost from \$2000-\$6000 per sample. If we assume similar rates for lake sampling, but
40 817 lower the cost as some aspects of lake sampling may be cheaper than stream sampling and multiply that
41 818 cost (estimated as \$1,000-\$4,000 US) by the total number of records of nutrient or chemical samples in
42 819 LAGOS-NE (n=589,909), then the combined estimate to collect the water quality data found in LAGOS-
43 820 NE is in the range of \$0.5 – 2.4 billion US. It cost us between 0.10 - 0.40 % of the cost to sample the data
44 821 in the first place to harmonize these half a million records, and to build an ecological-context database for
45 822 them. This relatively small investment in preserving, documenting, and harmonizing these valuable
46 823 datasets creates the needed infrastructure for new broad-scaled research, management, education, and
47 824 outreach uses.
48 825
49 826

50 825 Strategies for broad-scale data-integration efforts:

51 826 One challenge is to prioritize research areas and to identify the types of datasets that may benefit
52 827 from a similar type of integration. State, federal, tribal, and citizen-science water quality datasets were an
53 828 excellent source of quality data for integration and conducting broad-scaled research on aquatic systems.
54 829 There are likely other such data sources that would benefit from being integrated as we have done here.
55 830 We recommend the following strategies to make the best use of future data integration efforts.
56 831
57 832
58 833
59
60
61
62
63
64
65

1
2
3
4 833 (1) The database integration effort should be driven by key underlying research questions or goals, and
5 834 grounded in a strong conceptual foundation of the important features to include. In our case, the principles
6 835 of landscape limnology [18,19,20,12] guided the development of LAGOS-NE which helped us to
7 836 prioritize geospatial and lake features for inclusion in the database because the addition of any data type
8 837 or dataset cost time and money.
9 838

10 839 (2) For databases with more than one major data type, it is very helpful to build the database in modular
11 840 form, each with its own versioning system, specific data integration methods, and quality control
12 841 procedures. This strategy was not a primary goal at the outset of our project, but, it emerged somewhat
13 842 organically through the life of the project. We now recognize the many benefits that the modularity brings
14 843 to the database, including making it much easier to be dynamic rather than static by providing a platform
15 844 for the addition of new data, new types of data, and new modules in the future (such as for biological
16 845 data, or data from high-frequency sensors).
17 846

18 847 (3) The entire process should be grounded in an open-science framework. Knowing that the database,
19 848 design, and methods were to be shared and made usable by future users influenced our decisions
20 849 throughout the process, and made documentation a high priority throughout. Although we are making the
21 850 full database available now, before this point, we supported open science by publishing subsets of
22 851 LAGOS-NE data that were used in individual publications (e.g., [49, 50]).
23 852

24 853 (4) Creation of LAGOS-NE required a strong focus on team science, and in particular the roles of and
25 854 incentives for early-career researchers in such efforts. This type of research cannot be conducted in a
26 855 single-investigator mode, but requires a highly collaborative and effective team-based model (e.g., [51,
27 856 52, 53]). We explicitly considered strategies for ensuring that early-career team members get credit for
28 857 their contributions [54], and we recommend providing team members with opportunities for leadership,
29 858 project management, personnel management, and intellectual growth. For example, they can be part of
30 859 major decisions and can lead smaller efforts throughout the project, as well as be given power to shape
31 860 team policies and practices. This integration of early-career researchers into the entire research team and
32 861 effort will give early-career professionals deep knowledge of the database, the procedures, as well as the
33 862 skills to conduct such work in the future.
34 863

35 864 (5) The decision how to disseminate the database documentation needs to be considered early in the
36 865 project. For example, database documentation papers are rare, especially in ecology, but are very
37 866 important. The documentation and procedural approaches for developing this large, integrated, and
38 867 heterogeneous database had to be disseminated through publication prior to making the database
39 868 available [17] and prior to publication of research results stemming from LAGOS-NE because methods
40 869 sections in journal articles are too short to include all the necessary documentation of such methods.
41 870 Other researchers may be discouraged by the very real consequence that publishing such products take
42 871 time and energy investments that may slow down production of research publications. However, such a
43 872 paper was instrumental in supporting later research articles that used LAGOS-NE. Therefore, we
44 873 recommend that this (and other) database documentation papers become a more standard type of paper to
45 874 describe the extensive methods involved and to supplement data papers. Such papers will facilitate the
46 875 use, extension, and translation of these databases well into the future, as well as foster future research on
47 876 broad-scale, complex, and societally-relevant environmental questions.
48 877

878 **Availability of supporting source code and requirements**

879 Project name: LAGOS-NE

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

881 Operating system(s): e.g. Platform independent

882 Programming language: R

1
2
3
4 883 Other requirements: R packages required (with associated versions): dplyr ($\geq 0.7.0$), rappdirs ($\geq 0.3.1$),
5 884 lazyeval (≥ 0.2), purrr ($\geq 0.2.2.2$), magrittr (≥ 1.5), sf, curl ($\geq 2.7.0$), stringr ($\geq 1.2.0$)
6 885 License: GPL
7
8

9 887 **Availability of supporting data**

10 888 The datasets supporting the results of this article are available in the Ecological Data Initiative repository,
11 889 including the following specific components:

- 12 890 • LAGOS-NE-LOCUS v1.01 [55]
- 13 891 • LAGOS-NE-LIMNO v1.087.1 [56]
- 14 892 • LAGOS-NE-GEO v1.05 [57]
- 15 893 • LAGOS-NE-GIS v1.0 [58]
- 16 894 • Snapshots of the R package in the LAGOS GitHub page are also available in the GigaScience
17 895 repository, GigaDB [59].
18 896

20 897 **DECLARATIONS**

21 898 **List of abbreviations**

22 899 LAGOS-NE – LAke multi-scaled GeOSpatial and temporal database for the 17 Northeastern and
23 900 Midwest U.S. states

24 901 GIS – Geographic Information System

25 902 US-EPA – United States Environmental Protection Agency

26 903 EML – Ecological Metadata Language

27 904 USGS – United States Geological Survey

28 905 NHD – National Hydrography Dataset

29 906 WBD – Watershed Boundary Dataset

30 907 IWS – Interlake Watershed

31 908 SRP – Soluble Reactive Phosphorus

32 909 TP – Total Phosphorus

33 910 TN – Total Nitrogen

34 911 DOC – Dissolved Organic Carbon

35 912 MAV – Maximum Allowable Value

36 913 TDN – Total Dissolved Nitrogen

37 914 IQR – Interquartile Range

38 915 LULC – Land Use Land Cover

39 916 CONN – Connectivity and abundance (lake, stream, and wetland)

40 917 CHAG – Climate, Hydrology, Atmospheric deposition of nitrogen and sulfur, and surficial Geology

41 918 HUC – Hydrologic Unit Code

42 919 CSI – Cross Scale Interactions
43 920

47 921 **Competing interesting**

48 922 The authors declare they have no competing interests
49 923

50 924 **Funding**

51 925 The creation of LAGOS-NE was supported by:

52 926 The National Science Foundation MacroSystems Biology Program in the Emerging Frontiers Division of
53 927 the Biological Sciences Directorate (EF-1065786, EF-1638679, EF-1065649, EF-1065818, EF-1638554)
54 928 and the USDA National Institute of Food and Agriculture, Hatch project 176820 to PAS. KEW thanks the
55 929 STRIVE Programme (2011-W-FS-7) from the Environmental Protection Agency, Ireland. SMC thanks
56 930 the NSF Division of Biological Infrastructure (1401954).
57 931
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

The water quality data that are incorporated into LAGOS-NE were originally funded by the following sources:

934 State of Maine; Michigan Agricultural Experiment Station; Fisheries Division, Michigan Department of
935 Natural Resources; New York State Division of Water Quality; Wisconsin Department of Natural
936 Resources; University of Wisconsin-Madison; State/Trust; Michigan State University Agriculture
937 Experimental Station Disciplinary Research Grant Program; US EPA; US EPA Section 106/319 Grants;
938 Tribal General Fund; U.S. Army Corps of Engineers Federal Lakes Operation & Maintenance Funds;
939 Aquatic Plant Management Society; Aquatic Ecosystem Restoration Foundation; Michigan State
940 University; Michigan State University Department of Fisheries and Wildlife; EPA Star Fellowship to
941 K.S.C. (U-915342-01-0); Andrew W. Mellon Foundation; Federal Aid in Sport Fish Restoration Program
942 (Grant F-69-P, Fish Management in Ohio) administered jointly by the U.S. Fish and Wildlife Service and
943 the Ohio Department of Natural Resources, Division of Wildlife; Iowa Department of Natural Resources
944 (Contract #ESD04HALFasch110155); Minnesota Pollution Control Agency; NSF-Division of
945 Environmental Biology; Ohio DNR Division of Wildlife; University of Rhode Island Watershed Watch;
946 NSF Kellogg Biological Station Long Term Ecological Research (LTER) Program, DEB 1027253; NSF
947 North Temperate Lakes LTER Program, DEB 1440297; Lac du Flambeau Band and Bureau of Indian
948 Affairs; Indiana Department of Environmental Management; Missouri Department of Natural Resources;
949 Clean Water Act Section 16; Michigan Department of Environmental Quality; Massachusetts Water
950 Supply Protection Trust; US EPA Clean Air Markets Division (LTM Network); US EPA Office of
951 Research and Development; New York City Department of Environmental Protection (NYSDEP); City of
952 New York; USGS Water Availability and Use Science Program (WAUSP); U.S. Geological Survey; New
953 York State Energy Research and Development Authority; National Institute of Food and Agriculture,
954 U.S. Department of Agriculture, Hatch Grant 1003732; the New York State Department of Environmental
955 Conservation; Lake Sunapee Protective Association; National Oceanic and Atmospheric Administration;
956 Gull Lake Quality Organization; Clean Michigan Initiative; NSF grant DEB-1455461.

Authors' contributions

960 Data for the database were contributed by: LCB, MB, KEB, MGB, MTB, SRC, JWC, KSC, MC, JDC,
961 JAD, JD, CTF, CSF, MJG, LTG, JDF, SKH, PCH, EH, CH, JRJ, KJH, LLJ, WWJ, JRJ, CMK, SAK, BL,
962 JAL, YL, NRL, JAL, LJM, WHM, KEBM, BPN, SJN, MLP, DCP, AIP, DMP, POR, DOR, KMR, LGR,
963 OS, NJS, PAS, NRS, EHS, JLS, JMT, TPT, MV, GW, KCW, KEW, JDW, and MKW. The idea to create
964 the database was conceived by PAS and KSC. PAS coordinated the different activities across team-
965 members to build LAGOS-NE. The database was designed by EGB, PNT, CG, and PAS; and, created and
966 managed by EGB. The following authored metadata for the individual water quality data sets using
967 information provided by the data providers: MTB, CKB, KSC, SMC, CEF, CTF, ENH, NRL, SKO, NKS,
968 PAS, EHS, and KEW. CEF prepared the integrated LAGOS-NE metadata, and developed the protocols
969 for authoring the EML metadata; and, CEF and CKB created EML metadata for the 87 water quality data
970 sets. SKO wrote the final variables definitions for the integrated metadata. CG helped to prepare the
971 needed metadata and documentation for loading the data in the data repository. Code for importing the
972 datasets into the database was written by EGB, STC, NRL, and SY. NJS and SBS performed geospatial
973 analyses and created the LAGOS-GIS Toolbox. The conceptual foundation for measuring freshwater
974 connectivity was led by CEF. SBS developed the methods to delineate lake watersheds. The quality
975 control methods development and analysis on LAGOS-NE_{LIMNO} were conducted by NRL; the quality
976 control of LAGOS-NE_{GIS} was led by CES and SMC, and conducted by CES, SMC, CEF, NKS, and
977 KEW. The quality control of LAGOS-NE_{LOCUS} was conducted by EGB. Many authors who were part of
978 the database integration team wrote the technical documentation; JFL served as editor of these technical
979 documents. Tables and figures were prepared by SMC, KBSK, JFL, NRL, ACP, NKS and PAS and
980 edited by many of the contributing authors. SKO and JJS wrote the LAGOS-NE R package. NJS prepared
981 the GIS data and its corresponding metadata. PAS coordinated the writing of the manuscript, and major

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

982 parts of the manuscript were written by: PAS, KSC, SMC, JFL, NRL, SKO, JJS, EHS, PNT, TW, and
983 SY. After the lead author, authors are listed alphabetically.

984
985 **Competing Interests**

986 The authors declare that they have no competing interests

987
988 **Acknowledgments**

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

1
2
3
4 1002
5 1003
6 1004
7 1005
8 1006
9 1007
10 1008
11 1009
12 1010
13 1011
14 1012
15 1013
16 1014
17 1015
18 1016
19 1017
20 1018
21 1019
22 1020
23 1021
24 1022
25 1023
26 1024
27 1025
28 1026
29 1027
30 1028
31 1029
32 1030
33 1031
34 1032
35 1033
36 1034
37 1035
38 1036
39 1037
40 1038
41 1039
42 1040
43 1041
44 1042
45 1043
46 1044
47 1045
48 1046
49 1047
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

References

1. Carpenter SR, Caraco NR, Correll DL, Howarth RW, Sharpley AN, Smith VH. Nonpoint pollution of surface waters with phosphorus and nitrogen. *Ecological Applications*. 1998; 8:559-568.
2. Jaworski NA, Howarth RW, Hetling LJ. Atmospheric Deposition of Nitrogen Oxides onto the Landscape Contributes to Coastal Eutrophication in the Northeast United States. *Environmental Science and Technology* 1997; 31:1995–2004.
3. Bennett EM, Carpenter SR, Caraco NF. Human Impact on Erodible Phosphorus and Eutrophication: A Global Perspective. *BioScience*. 2001; 51:227-234.
4. Schindler DW. Recent advances in the understanding and management of eutrophication. *Limnology and Oceanography*. 2006; 51: 356–363.
5. Taranu ZE and Gregory-Eaves I. Quantifying relationships among phosphorus, agriculture, and lake depth at an inter-regional scale. *Ecosystems*. 2008; 11: 715-725.
6. Filstrup CT, Wagner T, Soranno PA, Stanley EH, Stow CA, Webster KE, Downing JA. Regional variability among nonlinear chlorophyll-phosphorus relationships in lakes. *Limnology and Oceanography*. 2014; 59: 1691-1703. doi: 10.4319/lo.2014.59.5.1691.
7. McCrackin ML, Jones HP, Jones PC, Moreno-Mateos D. Recovery of lakes and coastal marine ecosystems from eutrophication: A global meta-analysis. *Limnology and Oceanography*. 2016; 62:507-518. doi: 10.1002/lno.10441.
8. Paerl HW, Otten TG, Joyner AR. Moving towards adaptive management of cyanotoxin-impaired water bodies. *Microbial Biotechnology*. 2016; 9:641-651. doi: 10.1111/1751-7915.12383.
9. Schindler DW, Carpenter SR, Chapra SC, Hecky RE, Orihel, DM. Reducing Phosphorus to Curb Lake Eutrophication is a Success. *Environmental Science & Technology*. 2016; 50:8923-8929. doi: 10.1021/acs.est.6b02204.
10. Fergus CE, Soranno PA, Cheruvilil KS, Bremigan MT. Multiscale landscape and wetland drivers of lake total phosphorus and water color. *Limnology and Oceanography*. 2011; 56:2127-2146. doi: 10.4319/lo.2011.56.6.2127.
11. Soranno PA, Cheruvilil KS, Bissell EG, Bremigan MT, Downing JA, Fergus CE, Filstrup CT, Henry EN, Lottig NR, Stanley EH, Stow CA, Tan P-N, Wagner T, Webster KE. Cross-scale interactions: quantifying multi-scaled cause–effect relationships in macrosystems. *Frontiers in Ecology and the Environment*. 2014; 12:65–73.
12. Read EK, Patil VP, Oliver SK, Hetherington AL, Brenttrup JA, Zwart JA, Winters KM, Corman JR, Nodine ER, Woolway RL, Dugan HA, Jaimes A, Santoso AB, Hong GS, Winslow LA, Hanson PC, Weathers KC. The importance of lake-specific characteristics for water quality across the continental United States. *Ecological Applications*. 2015; 35:943-955. doi: 10.1890/14-0935.1.
13. Smith VH, Dodds WK, Havens KE, Engstrom DR, Paerl HW, Moss B, Likens GE. Comment: Cultural eutrophication of natural lakes in the United States is real and widespread. *Limnology and Oceanography*. 2014; 59: 2217-2225.
14. McDonald CP, Lottig NR, Stoddard, JL, Herlihy AT, Lehmann S, Paulsen SG, Peck DV, Pollard AI, Stevenson RJ. Comment on Bachmann et al. (2013): A non-representative sample cannot describe the extent of cultural eutrophication of natural lakes in the United States. *Limnology and Oceanography*. 2014; 59: 2226-2230.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

15. Stoddard JL, Sickle JV, Herlihy AT, Brahney J, Paulsen S, Peck DV, Mitchell R, Pollard AI. Continental-Scale Increase in Lake and Stream Phosphorus: Are Oligotrophic Systems Disappearing in the United States? *Environmental Science and Technology*. 2016; 50(7):3409-3415. doi: 10.1021/acs.est.5b05950.

16. Sprague LA, Oelsner GP, Argue DM. Challenges with secondary use of multi-source water-quality data in the United States. *Water Research*. 2017; 100: 252-261.

17. Soranno PA, Bissell EG, Cheruvilil KS, Christel ST, Collins SM, Fergus CE, Filstrup CT, Lapierree JF, Lottig NR, Oliver SK, Scott CE, Smith NJ, Stopyak S, Yuan S, Bremigan MT, Downing JA, Gries C, Henry EN, Skaff NK, Stanley EH, Stow CA, Tan PN, Wagner T, Webster KE. Building a multi-scaled geospatial temporal ecology database from disparate data sources: fostering open science and data reuse. *GigaScience*. 2015; 4:28. doi: 10.1186/s13742-015-0067-4.

18. Magnuson JJ, Kratz, TK. Lakes in the landscape: approaches to regional limnology. *International Association of Theoretical and Applied Limnology*. 2000; 27: 74-87.

19. Wiens JA. Riverine landscapes: taking landscape ecology into the water. *Freshwater Biology*. 2002; 47L501-515. doi: 10.1046/j.1365-2427.2002.00887.x.

20. Soranno PA, Cheruvilil KS, Webster KE, Bremigan MT, Wagner T, Stow CA. Using landscape limnology to classify freshwater ecosystems for multi-ecosystem management and conservation. *BioScience*. 2010; 60:440–454.

21. 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, Kramer DB, Bremigan MT, Ligmann-Zielinska A. Evaluating the effects of upstream lakes and wetlands on lake phosphorus concentrations using a spatially-explicit model. *Landscape Ecology*. 2012, 27:1015-1030; doi: 10.1007/s10980-012-9762-z.

23. United States Geological Survey Watershed Boundary Dataset. <https://nhd.usgs.gov/wbd.html>. Downloaded in 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. 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, Wagner T, Winslow LA, Scott CE, Webster KE, Downing JA, Stanley EH. Prediction of lake depth across a 17-state region in the United States. *Inland Waters*, 2016 6:314-324; doi: 10.5268/IW-6.3957.

29. Fergus CE, Lapierre JF, Oliver S, Skaff N, Cheruvilil K, Soranno, P, Webster K, Scott C. The freshwater landscape: Lake, wetland, and stream abundance and connectivity at macroscales. *Ecosphere*. 2017; 8:e01911. doi: 10.1002/ecs2.1911

30. Wagner T, Soranno PA, Cheruvilil KS, Renwick WH, Webster KE, Vaux P, Abbitt RJF. Quantifying sample biases of inland lake sampling programs in relation to lake surface area and land use/cover. *Environmental Monitoring and Assessment*. 2007; 141:131-147. doi:10.1007/s10661-007-9883-z.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

31. Strock KE, Saros JE, Nelson SJ, Birkel SD, Kahl JS, McDowell WH. 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; 127:353–365.

32. Seaber PR, Kapinos FP, Knapp GL. Hydrologic unit maps: U.S. Geological Survey water-supply paper 2294. U.S.G.S. 1987. <http://water.usgs.gov/GIS/huc.html>.

33. Smith NJ, Soranno PA, Stopyak S. LAGOS-NE GIS Toolbox. GitHub repository. 2014. https://soranno.github.io/LAGOS_GIS_Toolbox/.

34. Yuan S, Tan PN, Cheruvilil KS, Collins SM, Soranno PA. Constrained Spectral Clustering for Regionalization: Exploring the Trade-off between Spatial Contiguity and Landscape Homogeneity. *Data Science and Advanced Analytics*. 2015. doi:10.1109/DSAA.2015.7344878.

35. Cheruvilil KS, Yuan S, Webster KE, Tan PN, Lapierre JF, Collins SM, Fergus CE, Scott C, Henry E, Soranno PA, Filstrup CT, Wagner T. Creating multi-themed ecological regions for macroscale ecology: Testing a flexible, repeatable and accessible clustering method. *Ecology and Evolution*. 2017. 7: 3046–3058. doi: 10.1002/ece3.2884

36. Helsel DR. *Statistics for censored environmental data using Minitab and R*, 2nd edition. John Wiley and Sons, New York. 2012. <http://www.wiley.com/WileyCDA/WileyTitle/productCd-0470479884.html>.

37. Yuan S, Tan PN, Cheruvilil KC, Fergus CE, Skaff NK, Soranno PA. Hash-Based Feature Learning for Incomplete Continuous-Valued Data. *Proceedings of the 2017 SIAM International Conference on Data Mining*. Houston, Texas. April 27-29, 2017.

38. Collins SM, Oliver SK, Lapierre JF, Stanley EH, Jones JR, Wagner T, Soranno PA. Lake nutrient stoichiometry is less predictable than nutrient concentrations at regional and sub-continental scales. *Ecological Applications*. 2017; 27:1529-1540. doi: 10.1002/eap.1545

39. Fergus CE, Finley AO, Soranno PA, Wagner T. Spatial variation in nutrient and water color effects on lake chlorophyll at macroscales. *PLoS ONE*. 2016; 11(10). doi: 10.1371/journal.pone.0164592

40. Wagner T, Fergus EC, Stow CA, Cheruvilil KS, Soranno PA. The statistical power to detect cross-scale interactions at macroscales. *Ecosphere*. 2016; 7:e01417. doi: 10.1002/ecs2.1417.

41. Sterner RW, Elser JJ. *The Biology of Elements from Molecules to the Biosphere*. Princeton University Press: Princeton and Oxford; 2002.

42. Rastetter, EB. Modeling coupled biogeochemical cycles. *Frontiers in Ecology and the Environment*. 2001; 9: 68-73.

43. Finzi AC, Austin AT, Cleland EE, Frey SD, Houlton BZ, Wallenstein MD. Responses and feedbacks of coupled biogeochemical cycles to climate change: examples from terrestrial ecosystems. *Frontiers in Ecology and the Environment*. 2011; 9:61-67.

44. Finlay JC, Small GE, Sterner RW. Human influences on nitrogen removal in lakes. *Science*. 2013; 342:247-250.

45. Oliver SK, Collins SM, Soranno PA, Wagner T, Stanley EH, Jones JR, Stow CA, Lottig NR. Unexpected stasis in a changing world: Lake nutrient and chlorophyll trends since 1990. *Global Change Biology*. 2017. doi: 10.1111/gcb.13810

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

46. Lottig NR, Wagner T, Henry EN, Cheruvilil KS, Webster KE, Downing JA, Stow CA. 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. 2017. LAGOS: R interface to the LAke multi-scaled GeOSpatial & temporal database. R package version 1.087.1. <https://github.com/cont-limno/LAGOS>

48. Cheruvilil KS, Soranno P. Developing nutrient criteria in Michigan lakes: Revision, update, and validation of the lake-specific model for establishing expected nutrient conditions in Michigan lakes. Report to the Michigan Department of Environmental Quality. FY2012 205(j). 2015.

49. Oliver SK, Soranno PA, Fergus CE, Wagner T, Winslow LA, Scott CE, Webster KE, Downing JA, Stanley EH. 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.

50. Fergus CE, Finley AO, Soranno PA, Wagner T. 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](https://doi.org/10.6073/pasta/0ebd2e4c0705706b77b359955bff44e1).

51. Cheruvilil KS, Soranno PA, Weathers KC, Hanson PC, Goring SJ, Filstrup CT, Read EK. Creating and maintaining high-performing collaborative research teams: the importance of diversity and interpersonal skills. *Frontiers in Ecology and the Environment*. 2014; 12:31–38.

52. Weathers KC, Hanson PC, Arzberger P, Brentrup J, Brookes JD, Carey CC, Gaiser E, Hamilton DP, Hong GS, Ibelings B, Istvánovics V, Jennings E, Kim B, Kratz T, Lin FP, Muraoka K, O'Reilly C, Piccolo C, Ryder E, Zhu G. The Global Lake Ecological Observatory Network (GLEON): The Evolution of grassroots network science. *Bulletin of Limnology and Oceanography*. 2013; 22:71-73.

53. 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.

54. Goring S, Weathers KC, Dodds WK, Soranno PA, Sweet LC, Cheruvilil KS, Kominoski JS, Rüegg J, Thorn AM, Utz RM. Improving the culture of interdisciplinary collaboration in ecology by expanding measures of success. *Frontiers in Ecology and the Environment*. 2014; 14: 39-47.

55. Soranno P., K. Cheruvilil. 2017. 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. <http://dx.doi.org/10.6073/pasta/940b25d022c695b440e1bdbc49fbb77b>. Dataset accessed 9/14/2017.

56. Soranno P., K. Cheruvilil. 2017. 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. <http://dx.doi.org/10.6073/pasta/b1b93ccf3354a7471b93eccc484d506>. Dataset accessed 9/14/2017.

57. Soranno P., K. Cheruvilil. 2017. 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.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

<http://dx.doi.org/10.6073/pasta/b88943d10c6c5c480d5230c8890b74a8>. Dataset accessed 9/14/2017.

58. Soranno P., K. Cheruvilil. 2017. 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. <http://dx.doi.org/10.6073/pasta/fb4f5687339bec467ce0ed1ea0b5f0ca>. Dataset accessed 9/14/2017.

59. Bacon LC, Beauchene M, Bednar KE et al.: Supporting data for "LAGOS-NE: A multi-scaled geospatial and temporal database of lake ecological context and water quality for thousands of U.S. lakes" GigaScience Database. 2017. <http://dx.doi.org/10.5524/100350>

Figure 1. Map of the study extent of LAGOS-NE. Map includes 17 states in the upper midwest and northeastern U.S. outlined in white and 51,101 lakes ≥ 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 [17].

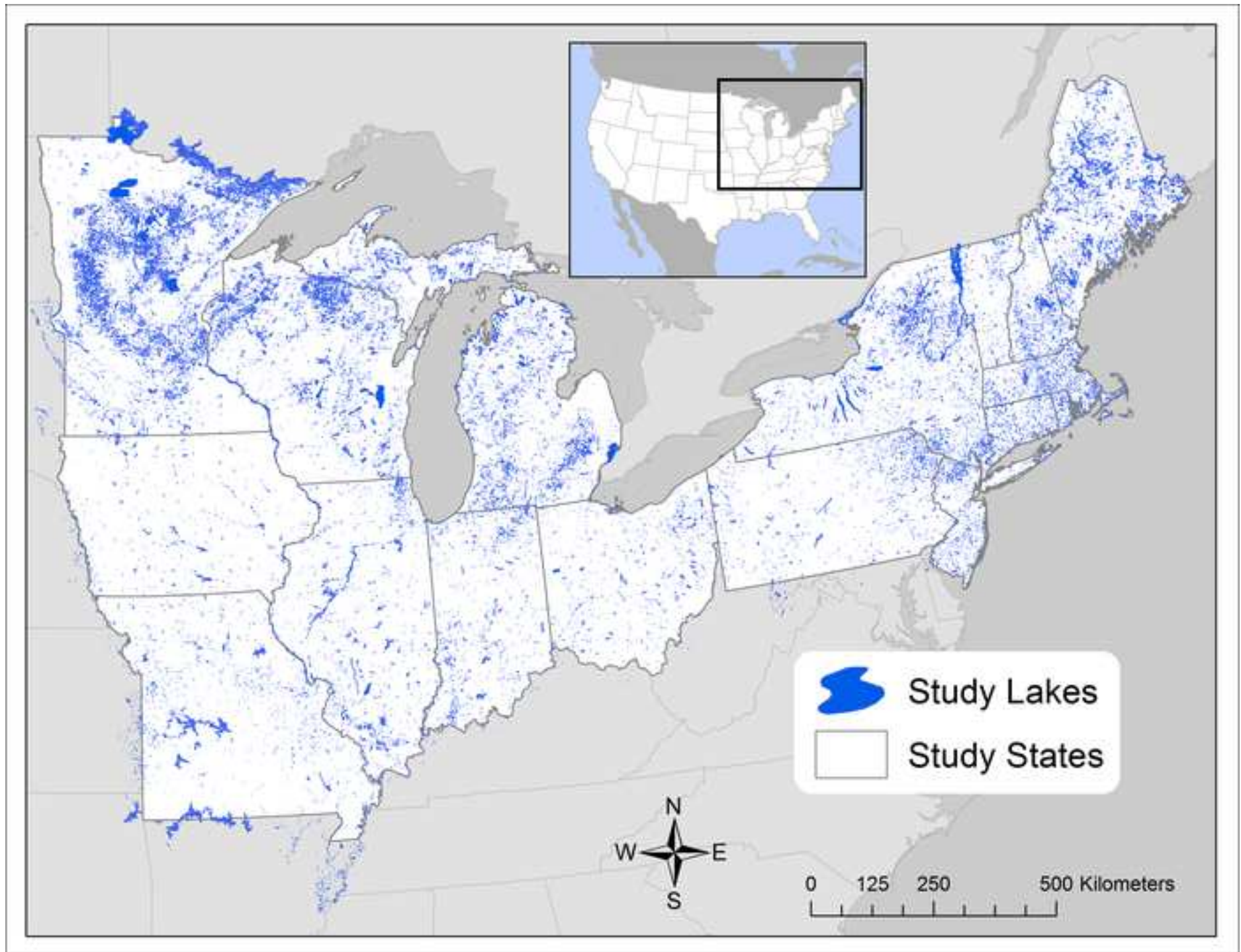
Figure 2. LAGOS-NE data modules and version numbers. The data modules and versions that are included in LAGOS-NE and are available with this paper include: LAGOS-NE_{GEO} v.1.05, LAGOS-NE_{LOCUS} v.1.01 (note, that in Soranno et al. [17], this module was called LAGOS-lakes), and LAGOS-NE_{LIMNO} v.1.087.1. We include descriptions of the type 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_{LOCUS} module (through the unique lake ID). P is phosphorus, N is nitrogen, C is carbon, S is sulfur, atm is atmospheric, NHD is the National Hydrography Dataset, IWS is the interlake watershed, WBD is the Watershed Boundary Dataset, EDU is Ecological Drainage Unit. Figure is modified from Figure 1 in Soranno et al. [17].

Figure 3. Examples lake watersheds (IWS) in LAGOS-NE. The watersheds are coded by hydrologic class to which its lake belongs. Data are from the LAGOS-NE_{GEO} v.1.01 data module and the GIS data coverages.

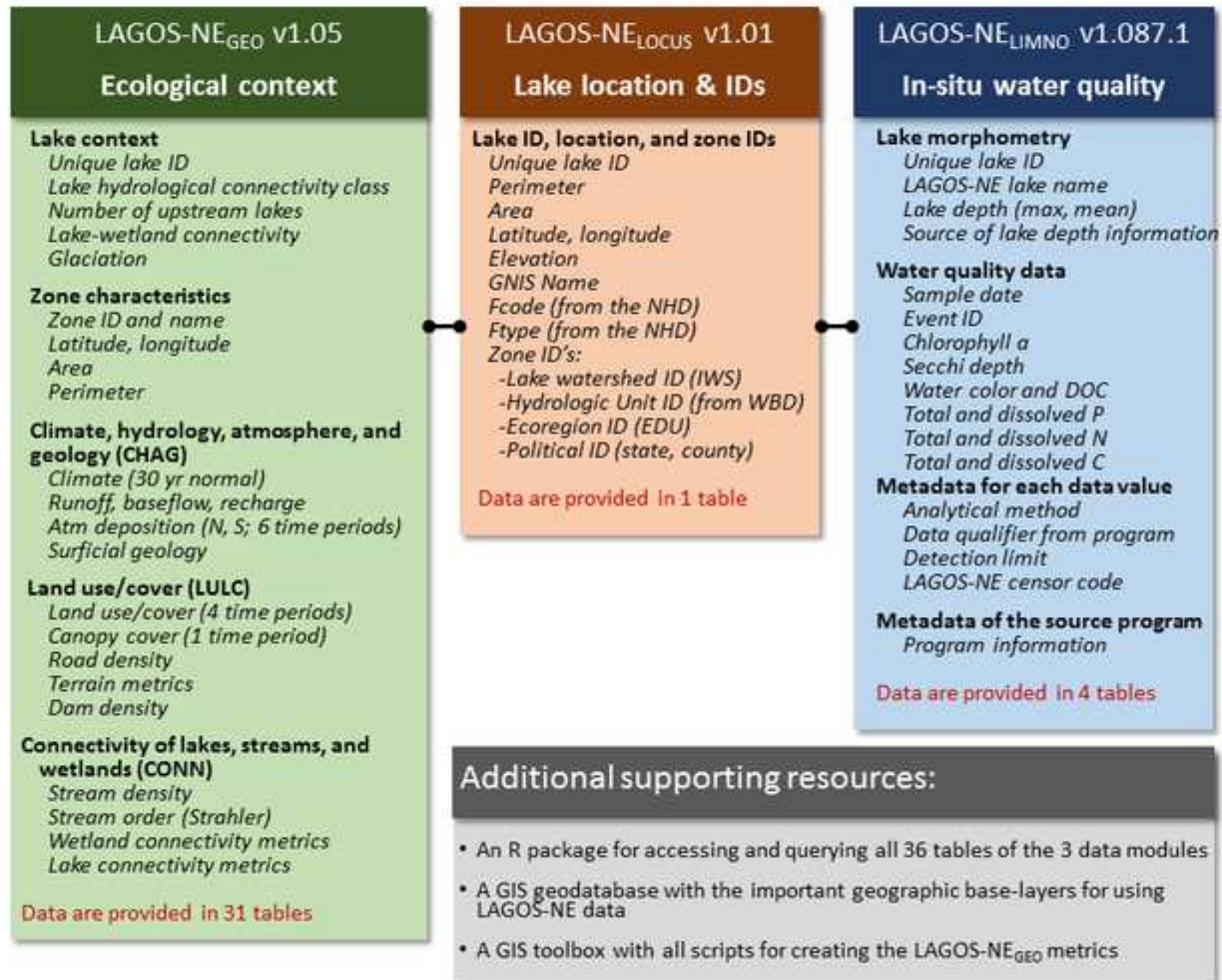
Figure 4. Percentage of lakes by lake area with water quality data. Percentage of census lakes in each lake size bin (top panel) compared to the percentage of census lakes for which there are limnological data for Secchi (second panel), chlorophyll *a* (third panel), and total phosphorus (TP; bottom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

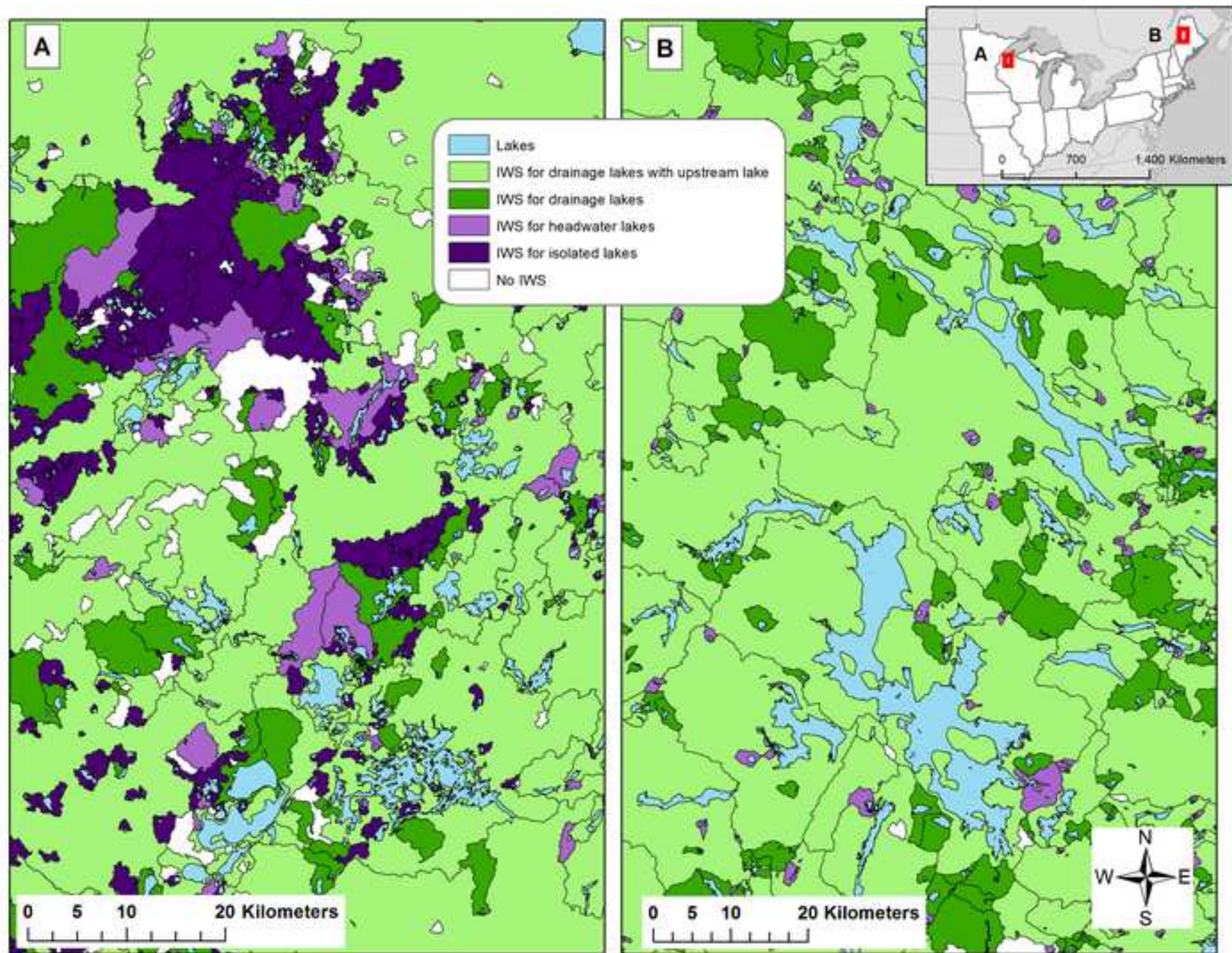
Figure 5. The number of years of water quality data by lake. The number of years for which at least one sample is taken during the summer stratified season (15 June to 15 September) for: Secchi depth in meters, total phosphorus in ug/L, total nitrogen in ug/L (includes both measured and calculated values),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

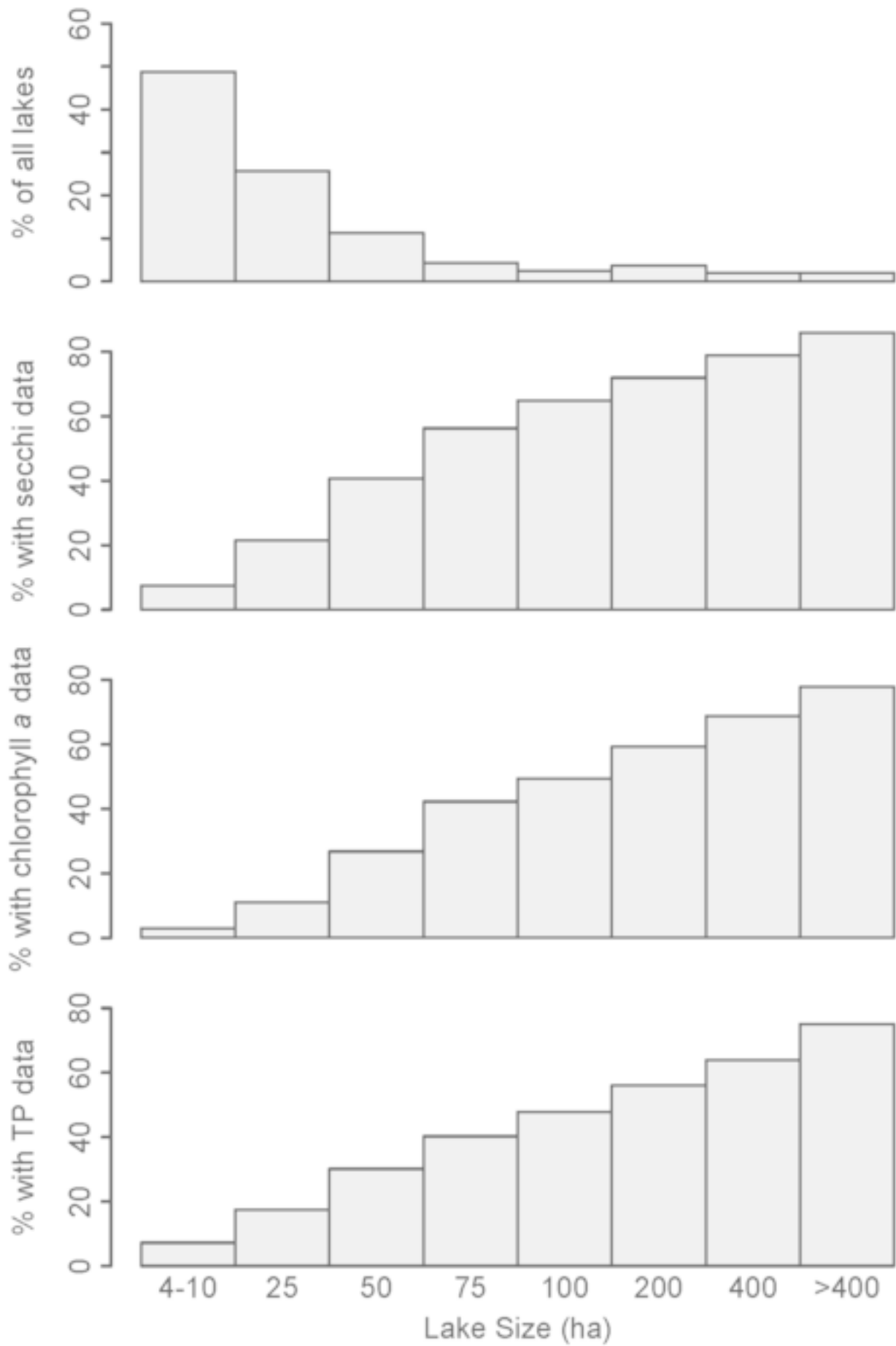
Figure 6. Example ecological context variables by spatial classification in LAGOS-NE. The top four panels are zoomed in to selected regions of Minnesota and Wisconsin so that the zone boundaries can be seen. The upper left panel shows stream density in each lake IWS, and the upper right panel shows the percent 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 percent of streamflow that is baseflow in each HUC8.

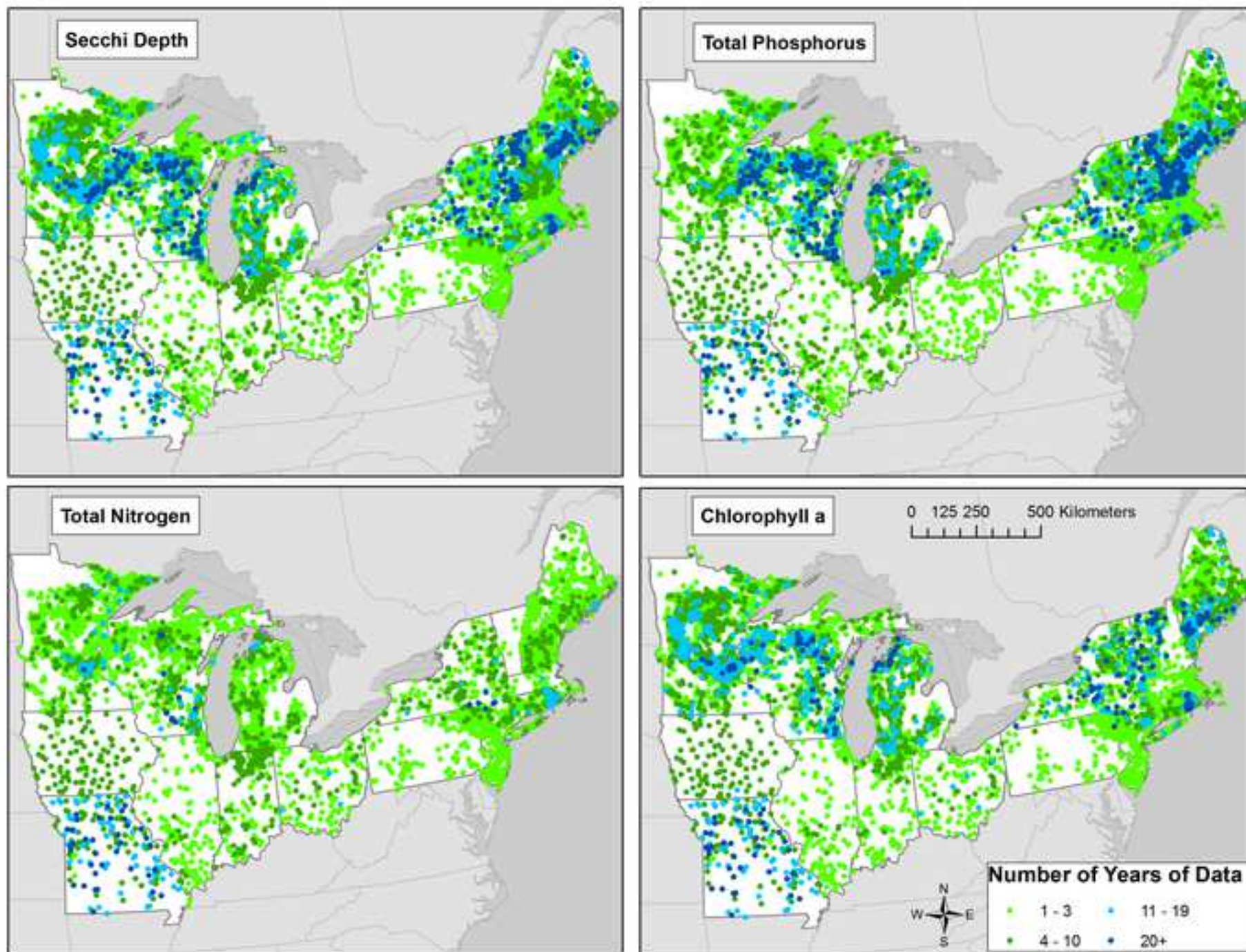


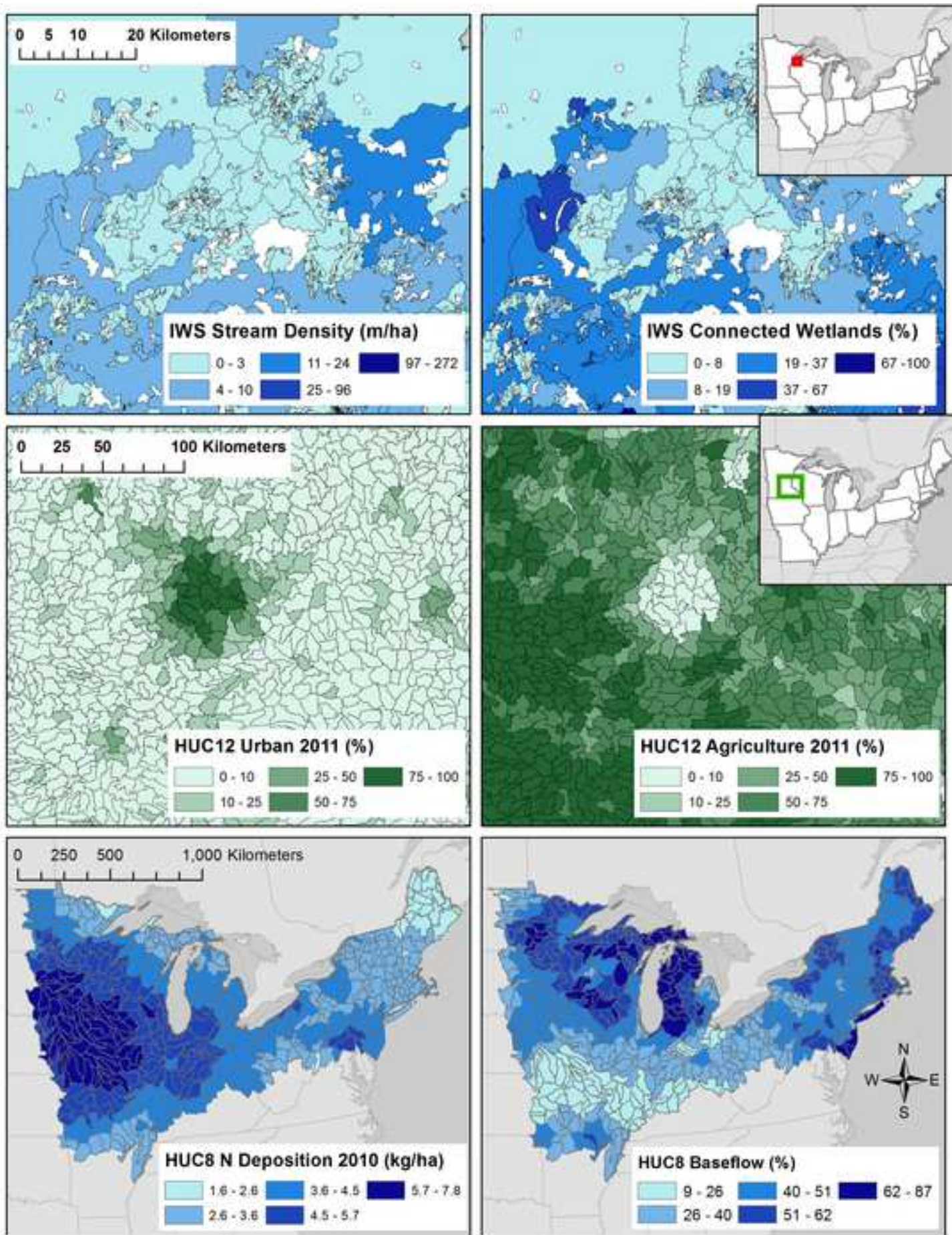
LAGOS-NE Modules













[Click here to access/download](#)

Supplementary Material

[Soranno_etal_2017_Additional file 1_8SEP17_v3.docx](#)

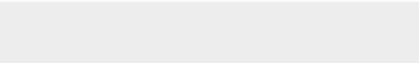





Click here to access/download

Supplementary Material

Soranno_etal_2017_Additional file 2_qaqc-
limno_v2.docx



Dr. Goodman, Editor-in-chief
GigaScience

May 16, 2017

Dear Dr. Goodman,

Please find enclosed our manuscript, *LAGOS-NE: A multi-scaled geospatial temporal database of lake ecological context and water quality for thousands of U.S. lakes* by Soranno et al., which we would like to submit for publication as a *Data Note* in GigaScience. We describe and make available a very large integrated, geospatial database for water quality. Given the size and scope of the database, the Data Note is likely longer than many such articles published in this journal. However, we wanted to include sufficient detail to correctly describe the database for future use of scientists. This is the database for which we described previously in a Review article published in GigaScience in 2015 (<https://gigascience.biomedcentral.com/articles/10.1186/s13742-015-0067-4>). For this Data Note, we are making the data available, and describing the data itself. Given the reception that the previous article received, we think the research community will find our database useful. Also, we appreciated the flexibility that GigaScience provided us in drafting our previous paper that describe the methods in detail, and having the data paper published in the same journal seems to make sense so that future data users can find both papers easily. We have decided to place the data in the Environmental Data Initiative data repository as described in our paper.

We confirm that this manuscript has not been published elsewhere and is not under consideration by another journal. All authors have approved the manuscript, agree with its submission to GigaScience, and have no competing interests.

Michigan State
University

College of
Agriculture and
Natural Resources

DEPARTMENT OF
FISHERIES AND
WILDLIFE

480 Wilson Road, Room 13
Natural Resources Bldg.
East Lansing, MI 48824

517/355-4478
FAX: 517/432-1699

We recommend the following reviewers:

Dr. Carly Strasser, DataCite Organization, Oakland, CA carlystrasser@gmail.com

Dr. Simon Goring, University of Wisconsin, Madison, goring@wisc.edu

Dr. Matthew Jones, University of California, Santa Barbara, jones@nceas.ucsb.edu

We look forward to hearing from you.

Sincerely,



Dr. Patricia A. Soranno, corresponding author
Professor
Michigan State University