

# A high-resolution inland surface water body dataset for the tundra and boreal forests of North America

Yijie Sui<sup>1</sup>, Min Feng<sup>\*1,2,3</sup>, Chunling Wang<sup>1,3</sup>, Xin Li<sup>1,2,3</sup>

<sup>1</sup>National Tibetan Plateau Data Center, Institute of Tibetan Plateau Research, Chinese Academy of Sciences, Beijing 100101, China

<sup>2</sup>CAS Center for Excellence in Tibetan Plateau Earth Sciences, Chinese Academy of Sciences, Beijing 100101, China

<sup>3</sup>University of Chinese Academy Sciences, Beijing 100049, China

*Correspondence to:* Feng M. (mfeng@itpcas.ac.cn)

**Abstract.** Inland surface waters are abundant in the tundra and boreal forests of North America, essential to environments and human societies but vulnerable to climate changes. These high-latitude water bodies differ greatly in their morphological and topological characteristics related to the formation, type, and vulnerability. In this paper, we present a water body dataset for the North American high latitudes (WBD-NAHL). Nearly 6.5 million water bodies were identified, with approximately 6 million (~90%) of them smaller than 0.1 km<sup>2</sup>. The dataset provides area and morphological attributes for every water body. During this study, we developed an automated approach for detecting surface water extent and identifying water bodies in the 10-m resolution Sentinel-2 multispectral satellite data to enhance the capability for delineating small water bodies and their morphological attributes. The approach was applied to the Sentinel-2 data acquired in 2019 to produce the water body dataset for the entire tundra and boreal forests in North America. The dataset provided a more complete representation of the region than existing regional datasets for North America, e.g., Permafrost Region Pond and Lake (PeRL). The total accuracy of the detected water extent by the WBD-NAHL dataset was 96.36% by comparing to interpreted data for locations randomly sampled across the region. Compared to the 30-m or coarser resolution water datasets, e.g., JRC GSW yearly water history, HydroLakes, and Global Lakes and Wetlands Database (GLWD), the WBD-NAHL provided an improved ability on delineating water bodies, and reported higher accuracies in the size, number, and perimeter attributes of water body by comparing to PeRL and interpreted regional dataset. This dataset is available from the National Tibetan Plateau/Third Pole Environment Data Center (TPDC, <http://data.tpdc.ac.cn>): DOI: 10.11888/Hydro.tpdc.271021 (Feng et al., 2020).

## 1 Introduction

Inland surface waters include various types of water bodies, including rivers and streams; large and small lakes; reservoirs; and ephemeral ponds. Inland surface water occupies only 2% of the global land surface (Pekel et al., 2016), but it plays a critical role in terrestrial ecosystems. Surface water distribution varies across the landscape. More than 55% of global surface waters are located in high latitudes in the Northern Hemisphere (> 44°N), and these northern high-latitude waters are generally small and densely clustered. The high latitudes have warmed faster than other regions, with annual surface temperatures increasing > 1.4° C over the past century (IPCC 2014). The temperature of the Arctic, in particular, has risen twice as fast as the average global temperature (Graversen et al., 2008; Johannessen et al., 2004; Pachauri and Reisinger, 2007; Serreze and Francis, 2006; Li et al., 2020). This change in climate is driving changes in terrestrial ecosystems in the Arctic as well. For example, increases in vegetation productivity have been observed across the northern high latitudes (Forkel et al., 2016). Meanwhile, high-latitude water bodies have started changing since the early 1970s (Carroll et al., 2011; Carroll and Loboda, 2017; Cooley et al., 2019; Smith et al., 2005; Fayne et al., 2020; Nitze et al., 2020). Although some changes are seasonal, and

37 therefore temporary, permanent changes have been reported, and small lakes in permafrost regions are found to be more  
38 vulnerable to permanent changes in water extent (Carroll and Loboda, 2017; Karlsson et al., 2014).

39 With observed rising temperatures (Biskaborn et al., 2019), permafrost thawing poses a threat to the stability of inland surface  
40 waters, especially in arctic lowland surface areas, where most of the water bodies could be thermokarst lakes (Jones et al.,  
41 2011; Olefeldt et al., 2016) and have strong interactions with permafrost in the regions. Thawing permafrost not only leads to  
42 the formation of lakes and ponds of various sizes, but also leads to the release of organic carbon in the form of carbon dioxide  
43 (CO<sub>2</sub>) and methane (CH<sub>4</sub>) (Serikova et al., 2019). Changes in lake formation may result in concomitant changes to the extent  
44 and connectivity of surface water bodies, which can greatly impact the sustainability of aquatic ecosystems.

45 The morphology of the water bodies could be shaped by the surrounding environment (Grosse et al., 2013; Laird et al., 2003;  
46 Schilder et al., 2013; Sharma et al., 2019; Carpenter, 1983; Higgins et al., 2021). Shoreline complexity affects lake ice  
47 formation (Sharma et al., 2019). Lake connectivity affects fish migration (Laske et al., 2019; McCullough et al., 2019), fish  
48 habitats, and aquatic assemblages (Napiórkowski et al., 2019; Jiang et al., 2021), water self-purification and accelerates water  
49 cycling (Glińska - Lewczuk, 2009; Vaideliene & Michailov, 2008; Xiong et al., 2017). The density of water bodies impacts  
50 fish density and biomass (Sandlund et al., 2016; van Zyll de Jong et al., 2017; King et al., 2021). The shape and distribution  
51 of water bodies reflect what led to the water body formation (Laurence C. Smith et al., 2007). Furthermore, information about  
52 lake area extent can improve arctic land surface modeling (Langer et al., 2016; van Huissteden et al., 2011). For these reasons,  
53 it is critical to quantify high-latitude surface water extent, as well as characterize related morphological and topological features,  
54 including size and shape.

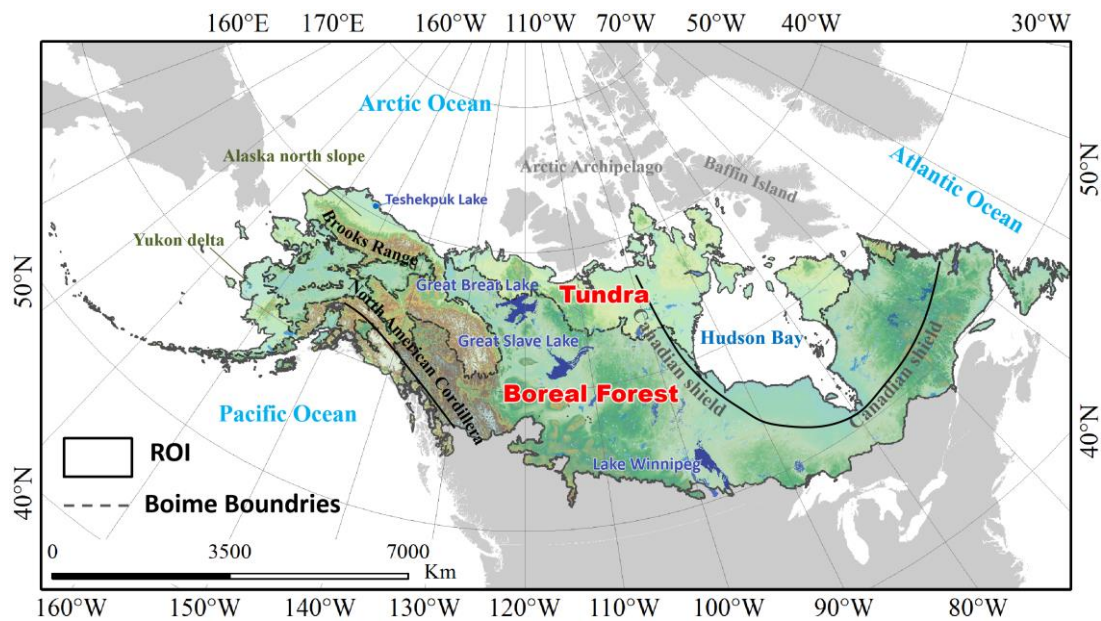
55 In the past, inland surface water was mapped at sub-hectare (i.e., 30 m) resolution using satellite data (Feng et al., 2015; Pekel  
56 et al., 2016; Pickens et al., 2020), and these data provided unprecedented information about the global extent of inland waters,  
57 including their spatial distribution and temporal changes. These datasets provide data that delineates the extent of large and  
58 moderate sizes of water bodies but underrepresent or fail to include the large number of small water bodies. Coarse-resolution  
59 datasets also lead to underrepresentation in delineating complex shorelines and the shapes of surface water bodies, making it  
60 difficult to derive their morphological and topological attributes. Existing datasets containing information that describe water  
61 body shapes, such as the Global Lakes and Wetlands Database (GLWD) (Lehner and Döll, 2004) and HydroLAKES (Messenger  
62 et al., 2016) are limited to water bodies larger than 0.1 km<sup>2</sup>. In spite of these limitations, these datasets provide valuable  
63 information for improving the precision of mapping inland waters. Detecting the extent of inland surface water at finer spatial  
64 scale boosts our ability to map small water bodies and improves the precision of delineating the shorelines of water bodies.  
65 This analysis then allows us to derive an inventory dataset of water bodies along with their morphological and topological  
66 attributes. The information allows scientists to analyze a water body as an object instead of a cluster of pixels, advancing our  
67 analysis and understanding of the water bodies' size, shoreline complexity, ecological effects, hydrological function, and  
68 vulnerability to natural and anthropogenic changes.

69 In this paper ,we present a higher resolution water body dataset for the North American high latitudes (WBD-NAHL). The  
70 dataset was derived by identifying the extent of inland waters using 10-m resolution Sentinel-2 multispectral data. The dataset  
71 provides the spatial extent and morphological attributes for each identified water body. It is the first inland water inventory  
72 dataset derived at this landscape scale with the capability of delineating inland surface waters as small as 0.001 km<sup>2</sup>.

## 73 **2 Spatial extent**

74 The WBD-NAHL dataset covers all tundra and boreal forest biomes in North America (Figure 1), with the exception of the  
75 Arctic Archipelago and Baffin Island due to their long time of snow or ice covering over water bodies. The topography of the

76 tundra and boreal forest in North America is extremely diverse, varying from mountains and rolling hills to plateaus and flat  
 77 coastal plains. The mountains of the North American Cordillera are covered by numerous mountain glaciers and also a large  
 78 number of glacial lakes. A large number of thermokarst lakes were found in lowland tundra areas, e.g., the Yukon Delta and  
 79 the Alaska North Slope (Olefelt et al., 2016). The vast Canadian Shield also has a high density of lakes. The climate of this  
 80 study region is characterized by long, cold winters and short, cool summers. The plants in the northern tundra include lichen,  
 81 moss, grass, sedge, and shrub. The southern boreal forest is dominated by evergreen forests (Ritter, 2006). Lakes are widely  
 82 distributed in the study region and approximately 36% of the land surface is covered by water (Messenger et al., 2016). The  
 83 number of lakes in this region accounts for 50% of the global lakes, and the area of lakes accounts for 30% of the global lakes  
 84 in the region, indicating the region to be one of the richest areas of surface water bodies (Messenger et al., 2016). Various types  
 85 of lakes, including organic, fluvial, meteorite, volcanogenic, and anthropogenic lakes, are distributed in the study region and  
 86 feature very different sizes and shapes (Dranga et al., 2017).



87  
 88 **Figure 1: The extent of the study area, including the tundra and boreal biomes, in the North Americas continent, excluding the**  
 89 **Arctic Archipelago and Baffin Island.**

90 **3 Data**

91 **3.1 Sentinel-2 A/B multi-spectral images**

92 Sentinel-2 multi-spectral images were used to delineate surface water bodies in this study. Sentinel-2 A/B provides a short  
 93 revisit cycle (2-3 days) in the high latitudes, which is critical for detecting surface water during the short, snow-free season in  
 94 the region. Sentinel-2 images were obtained using the United States Geological Survey (USGS) EarthExplorer client/server  
 95 interface (<https://earthexplorer.usgs.gov/>, last access: 7 April 2021).

96 Each Sentinel-2 image consists of 13 multispectral bands, including four bands at 10-m resolution, six bands at 20-m resolution,  
 97 and three others at 60-m resolution. Sentinel-2 data are distributed as collections representing different processing levels. We  
 98 selected the Sentinel-2 Collection 2 data, which provides spectral bands of surface reflectance after atmospheric corrections.  
 99 The 10-m Sentinel-2 bands were used for water detection to maximize spatial precision for delineating small water bodies.  
 100 The 20-m Sentinel-2 bands were resampled to 10-m resolution to match the higher resolution bands. The “s2cloudless”  
 101 (<https://github.com/sentinel-hub/sentinel2-cloud-detector>, last access: 7 April 2021) was applied to identify cloud-  
 102 contaminated pixels, generating a probability of cloud and cirrus detection. This module includes a model generated by a

103 Convolutional Neural Networks (CNN) trained with 6.4 million manually labeled samples. This model was validated to have  
104 99% accuracy for identifying clouds and 84% accuracy for identifying cirrus in Sentinel-2 images (Zupanc, 2020).

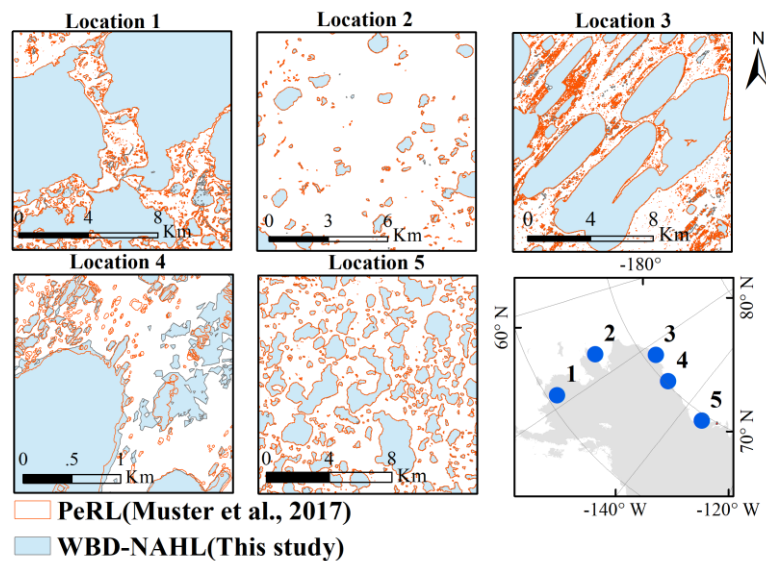
### 105 3.2 Joint Research Centre (JRC) yearly water dataset

106 The JRC yearly water dataset (JRC GSW Yearly Water Classification History, v1.2, <https://global-surface-water.appspot.com/>)  
107 (Pekel et al., 2016) provides a delineation of permanent water, non-water, and seasonal water for global inland surface waters.  
108 The dataset was produced using long-term Landsat images, including Landsat TM, ETM+, and OLI images acquired from  
109 1984 to 2019. Permanent water in the dataset was identified as water cover throughout the entire year, and seasonal water is  
110 identified based on occurrence during a single year.

111 The JRC yearly water dataset provides a reasonably accurate delineation of water distribution for the period 1984-2019, but  
112 its precision is limited by the 30-m spatial resolution of Landsat data. The dataset's accuracy at high latitudes is affected by  
113 the relatively poor return cycle of Landsat (16 days), cloudiness, and long periods of snow and ice in the region each year. The  
114 JRC dataset was used as a reference to overcome these limitations and improve our ability to identify and monitor inland  
115 surface water bodies, particularly small water bodies. The permanent water class in the JRC dataset was used in this analysis,  
116 while the seasonal water was excluded due to its reportedly low accuracy (Meyer et al., 2020). The maximum extent of  
117 permanent water bodies for the time period 1984-2019 were processed to fill gaps in individual years, which were then used  
118 as the reference in this study.

### 119 3.3 Permafrost Region Pond and Lake (PeRL)

120 The Permafrost Region Pond and Lake (PeRL) dataset was produced through a circum-Arctic effort to map ponds and lakes  
121 from modern (2002–2013) high-resolution aerial and satellite imagery with a resolution of 5 m or finer, including imagery  
122 from GeoEye, QuickBird, WorldView-1/2, the KOMPSAT-2, and TerraSAR-X. The PeRL dataset includes 69 small maps  
123 representing a wide range of environmental conditions in tundra and boreal biomes (Muster et al., 2017). There are 14 maps  
124 mainly distributed in five regions of North America (Figure 2). Because of the high-resolution data, the PeRL dataset is able  
125 to delineate water bodies as small as  $10^{-7}$  km<sup>2</sup>, which is valuable for validating satellite-derived water datasets for regions  
126 dominated by small water bodies.

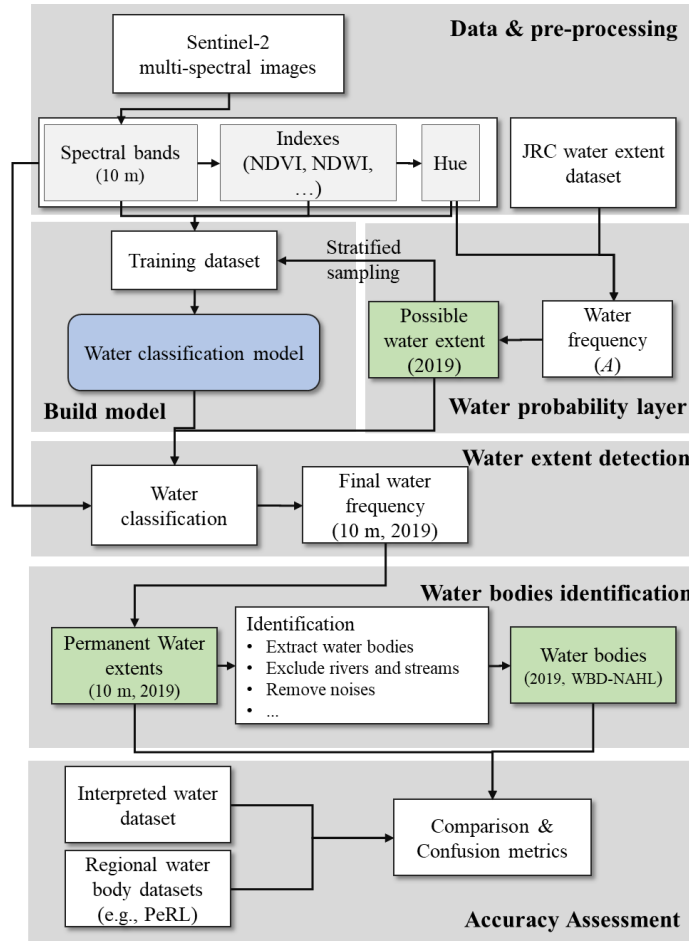


127

128 **Figure 2: Water bodies identified in the WBD-NAHL (this study) and PeRL datasets (Muster et al., 2017), and the locations (blue**  
129 **dots) of the PeRL maps for the study region.**

130 **4 Methods**

131 The 10-m resolution Sentinel-2 A/B multispectral data are the primary source used to identify small water bodies. An approach  
132 was developed to produce a water probability layer for 2019 by combining the water-sensitive indexes derived from the  
133 Sentinel-2 bands and the 30-m resolution JRC water dataset (section 4.1). A machine learning model was trained to retrieve  
134 water extent from the Sentinel-2 images from possible water extent restricted by the water probability layer (section 4.2)  
135 (Figure 3). Water bodies were finally identified from the water extent using an object-based algorithm to produce the final  
136 water body inventory (section 4.3).



137  
138 Figure 3: Flowchart for processing water extent and identifying water bodies.

139 **4.1 Water probability layer**

140 A water probability layer was derived to represent the likelihood of a pixel to correspond to permanent water during the  
141 summer of 2019. The 10-m resolution water-sensitive indexes calculated from the Sentinel-2 multispectral bands were used  
142 as the main input. The other reference water dataset (e.g., the JRC water dataset) was adopted as a supplemental input and  
143 fused with the main input to produce the water probability estimate at each 10-m resolution pixel.

144 To reduce effects of snow cover, Sentinel-2 A/B images acquired between June and September 2019 were selected to represent  
145 the relatively snow-free season in North American tundra and boreal biomes. The pixels in each Sentinel-2 image with an  
146 estimated cloud probability higher than 65% were excluded to avoid the effects of cloud contamination.

147 During pre-processing, multiple water-sensitive indexes were derived from each Sentinel-2 image to enhance the ability to  
148 detect water (Figure 3). To maximize the ability to separate water from non-water, especially vegetated land, three indexes  
149 were calculated to represent water and vegetation in each image: Normalized-Difference Water Index (NDWI) (McFeeters,

150 1996), Normalized Difference Vegetation Index (NDVI) (Carlson and Ripley, 1997), and Modified Normalized-Difference  
151 Water Index (MNDWI) (Xu, 2006). The three indexes were calculated as follows.

$$152 \quad NDWI = (B_{green} - B_{nir}) / (B_{green} + B_{nir}), \quad (1)$$

$$153 \quad NDVI = (B_{nir} - B_{red}) / (B_{nir} + B_{red}), \quad (2)$$

$$154 \quad MNDWI = (B_{green} - B_{swir}) / (B_{green} + B_{swir}), \quad (3)$$

155 Where  $B_{green}$ ,  $B_{red}$ ,  $B_{nir}$ , and  $B_{swir}$  are green (band #3), red (band #4), near-infrared (band #8), and short-wave infrared  
156 (band #11), respectively. These bands have 10-m resolution except  $B_{swir}$ , which has 20-m resolution and was pan-sharpened  
157 using the À Trous Wavelet Transform (ATWT) algorithm as recommended by Du et al., (2016). An HSV color space  
158 conversion was used to combine the three indexes and produce a final index for identifying water. The HSV (hue-saturation-  
159 value) color space conversion is a non-trigonometric pair of transformations from a linear red-green-blue (RGB) color space  
160 to a perceived color space (Danielson and Gesch, 2011). This method converts the three input bands into hue (color), saturation,  
161 and value components. The three indexes (NDWI, MNDWI, and NDVI) were scaled by 255, converted to a byte value type,  
162 combined into the RGB color space, and then converted to the HSV color space to derive a comprehensive index for identifying  
163 water.

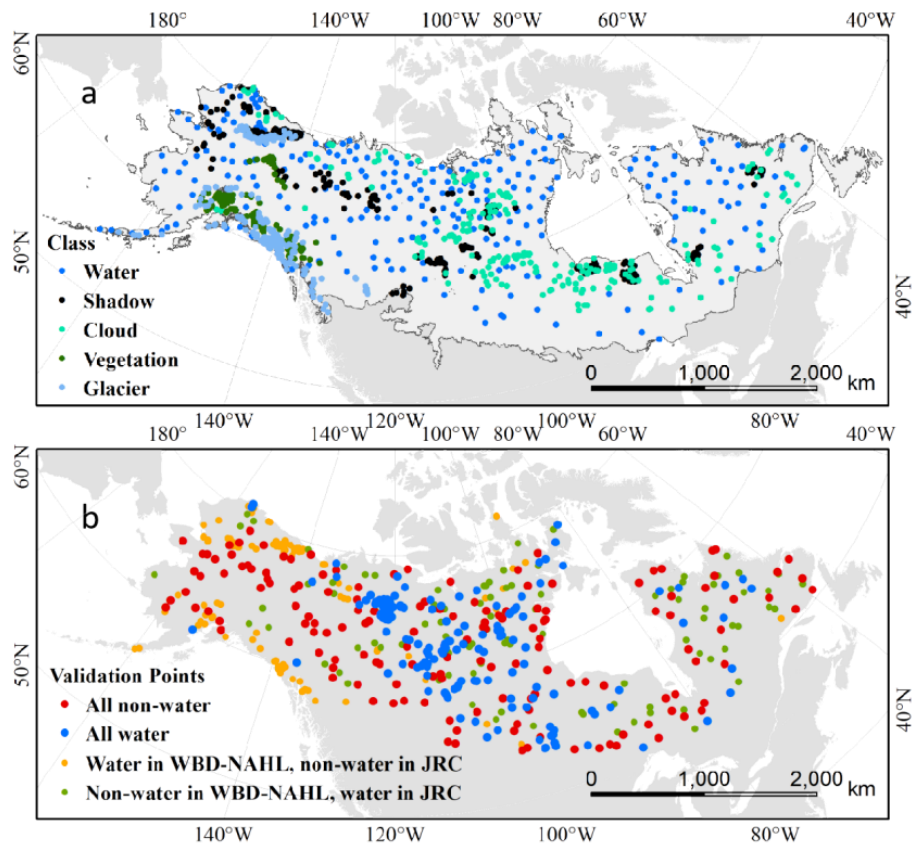
164 Once the hue has been identified, an experimental threshold of  $< 0.45$  was applied to identify the water pixels. The same  
165 procedure was applied to derive temporal water extents from all selected Sentinel-2 images. All the water extents were then  
166 combined to calculate the water frequency ( $A_s$ ) for the year. Potential water extent was then derived from the calculated water  
167 frequency data. The existing JRC water dataset provided complementary information for estimating possible water extent. The  
168 JRC permanent water records were resampled to 10-m resolution using the nearest neighbor algorithm and combined with the  
169 Sentinel-2-derived water frequency dataset using a weighted linear combination:

$$170 \quad A = W_s \cdot A_s + (1 - W_s) \cdot A_j, \quad (4)$$

171 where,  $A$  is the updated water frequency,  $W_s$  is the weight for the Sentinel-2-derived water frequency ( $A_s$ ) and set to 0.85 to  
172 ensure that the 10-m measurements were the main input for the final water probability estimate. However,  $W_s$  was decreased  
173 to 0.65 in high elevations pixels (elevation  $> 1$  km) to reduce the effect of snow and ice on the Sentinel-2-derived hue over  
174 mountains.  $A_j$  is the JRC permanent water record, which was set to 1.0 for permanent water and to 0.0 for others. The final,  
175 combined possible water extent was identified when  $A > 0.5$ .

## 176 4.2 Water extent detection

177 Although the possible water extent estimated the likelihood of a pixel to correspond to water, confusion with shadow, ice, or  
178 cloud contamination in area with complex environments is still possible due to the limitations of water indexes with similar  
179 spectra (Isikdogan et al., 2017). A random forest model was trained with points collected through visual interpretations to  
180 further detect water within the areas indicated as possible water. To ensure the representation of water and other land covers  
181 that can easily be confused as water, five strata were introduced, i.e., water, glacier, mountain, vegetation, and cloud. Then,  
182 250 points were randomly selected in each stratum, for a total of 1,250 points (Figure 4a).



183

184 **Figure 4: Training samples for random forest model building (a) and points identified for validating the accuracy of the detected**  
 185 **water extent (b).**

186 The five strata were established using reference datasets or customized rules. The glacier stratum was identified using the  
 187 Global Land Ice Measurements from Space (GLIMS) dataset of 2017 (<http://www.glims.org/>, last access: 7 April 2021), which  
 188 was a dataset of global glacier outlines including glacier area, geometry, surface velocity, and snow line elevation and was  
 189 produced from the Advanced Spaceborne Thermal Emission and Reflection Radiometer (ASTER) and the Landsat Enhanced  
 190 Thematic Mapper Plus (ETM+), as well as historical information derived from maps and aerial photographs. Vegetation was  
 191 identified as areas with a positive mean NDVI value calculated from the June-September Sentinel-2 images. The cloud stratum  
 192 was identified as having at least 20% of mean cloud probability calculated from the selected Sentinel-2 images. The mountain  
 193 shadow stratum was identified as any elevation higher than 1-km and slope greater than or equal to 3°. The water stratum was  
 194 identified as the remaining area of possible water extent.

195 The selected points were interpreted by the team to provide training data. Although we only used Sentinel-2 images from June  
 196 to September 2019, points were matched with a randomly selected image at the location during the time period, providing  
 197 representation for possible temporal variation. Each point was visually labeled by an interpreter after examining the image.  
 198 Metrics for visible bands (red, green, and blue), NDWI, MDWI, NDVI, and hue were derived from each image to provide  
 199 attributes for the point. These attributes were pooled to produce training data for building the machine learning model.

200 The scikit-learn Random Forest algorithm (Breiman, 2001) was adopted to build the model for surface water detection. This  
 201 model was applied to the selected Sentinel-2 images to detect surface water pixels. The results were compiled temporally to  
 202 produce a water frequency layer (*f*).

203 In this study, terrain shadows in the water frequency layer were removed with a terrain mask derived from the Global Multi-  
 204 resolution Terrain Elevation Data (GMTED) (Danielson and Gesch, 2011). The mask was where the slope was greater than or  
 205 equal to 7° and the elevation was over 1 km. The elevation threshold was used to minimize the impact of the slope threshold

206 on rivers in lowlands. The method using slope to identify terrain shadows was verified to be more effective than using hill-  
207 shade (Carroll and Loboda, 2017).

### 208 4.3 Water bodies identification

209 Permanent water pixels were identified from the resulting water frequency layer ( $f$ ) as being those pixels with at least 50%  
210 occurrence between June and September. The resulting water pixels were then converted to vector polygons using the “Raster  
211 to Polygon” tool in ESRI ArcMap 10.2. These water polygons provided the preliminary surface water body records.

212 An array of geometry metrics was calculated for each water body polygon using ArcMap in the  
213 Canada\_Lambert\_Conformal\_Conic projection (datum D\_North\_American\_1983 and Spheroid GRS80). These metrics  
214 include area, perimeter, and a shape index ( $SI$ ), which estimates the complexity of a water body polygon. The  $SI$  was calculated  
215 as:

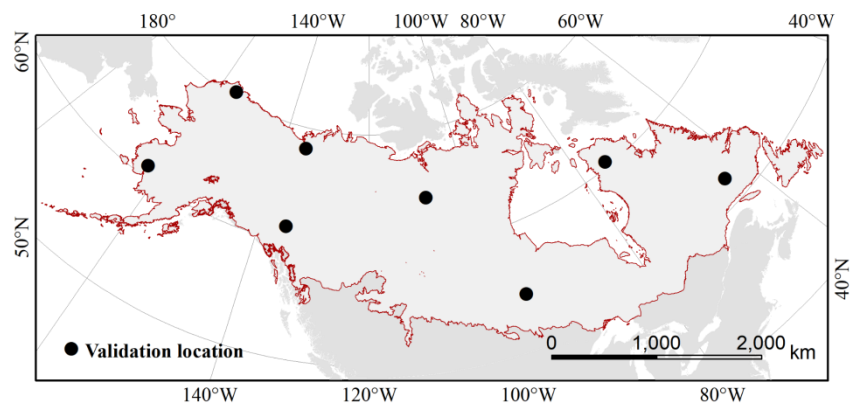
$$216 SI = P_{water_i} / P_{circle_i}, \quad (5)$$

217 where  $P_{water_i}$  is the perimeter of the water body  $i$ ,  $P_{circle_i}$  is the perimeter of a circle that has the same area as water body  $i$ .  $SI$   
218 equals 1 when a polygon is a perfect circle and greater than 1 when the polygon has a complex irregular shape.

219 The derived water body morphological metrics (i.e., the  $SI$  and area) and the HydroRIVERS were used to identify rivers and  
220 streams in the WBD-NAHL water bodies. Rivers and streams tend to have long, narrow, and linear shapes. We applied area  
221 thresholds  $> 5 \text{ km}^2$  and  $SI > 10$  in combination with visual examination to exclude large rivers and streams in WBD-NAHL.  
222 Considering the extreme difficulties in distinguishing small rivers and streams, water bodies that could possibly be rivers and  
223 streams were further identified by selecting long and linear water bodies ( $SI > 3$ ) located close to the rivers and streams ( $< 100$   
224 m), as indicated by HydroRIVERS.

### 225 4.4 Quality assessment

226 The accuracy and uncertainty of WBD-NAHL were assessed at two levels, i.e., pixel water extent and derived water bodies,  
227 to provide a comprehensive evaluation of the dataset. We randomly selected eight square blocks with a size of 10 km by 10  
228 km in the North American tundra and boreal region (Figure 5). The selected blocks were visually interpreted by the team to  
229 identify all the water bodies within each using a high-resolution Google Earth image as reference for interpretation. Water  
230 bodies records from the PeRL were compared to the WBD-NAHL water bodies to assess the number of water bodies and  
231 spatial area of each. The interpreted dataset was also compared to the JRC-derived water body records for 2019 to assess its  
232 accuracy in terms of representing water bodies. The JRC dataset provides water/nonwater map at the 30-m resolution pixels,  
233 representing the distribution of water extent, but no information in the spatial relationship between pixels and water bodies  
234 were provided, and we derived water bodies records from the JRC dataset using the same algorithm described in section 4.1.





236 **Figure 5: Locations of the five regions selected and interpreted for assessing the accuracy of the indicators of water bodies.**

237 The 14 regional PeRL maps were compared to the WBD-NAHL water bodies. Although the PeRL maps were produced from  
238 high-resolution images acquired in 2002-2013, the maps show little temporal changes when comparing to the WBD-NAHL  
239 dataset in the extents of the maps (Figure 2), and these maps were adopted as references for evaluating the WBD-NAHL water  
240 bodies. The PeRL maps were produced from images with 5 m resolution or finer, we excluded all water bodies in PeRL smaller  
241 than 0.0003 km<sup>2</sup> to ensure comparability to the scale of the WBD-NAHL dataset.

242 The water extent derived from the Sentinel-2 images were assessed by manually comparing specific points between the WBD-  
243 NAHL dataset and the JRC surface water dataset. The points were collected using a stratified random sampling across the  
244 entire study region. To achieve higher sampling performance, the outcomes were divided into four strata that represent pixels  
245 that were agreed as water, disagreed as water, agreed as non-water, and disagreed as non-water. In each of the strata, 400 points  
246 were randomly selected from the dataset and manually assessed by examining the same point in the latest Google Earth image.  
247 (Figure 4b) The results from the 1600 points were compared to the derived water extent. The confusion matrix was calculated  
248 from the results.

249 The sampling weights were included in the calculation of the metrics as following:

$$250 W_s = A_s/A_{all}, \quad (6)$$

251 where  $A_s$  is the area of stratum  $s$ , and  $A_{all}$  is the total area of the region.

252 Equations of the confusion metrics with weights:

$$253 OA = \sum_s^4 W_s * OA_s, \quad (7)$$

$$254 UA = \sum_s^4 W_s * UA_s, \quad (8)$$

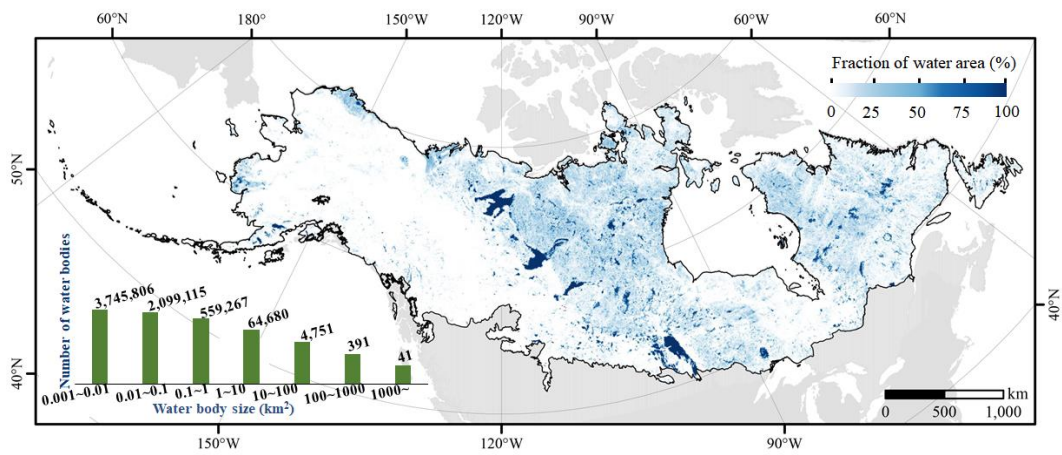
$$255 PA = \sum_s^4 W_s * PA_s, \quad (9)$$

256 where  $OA$ ,  $UA$ , and  $PA$  are the overall accuracy, user's accuracy and producer's accuracy of the entire dataset,  $OA_s$ ,  $UA_s$ , and  
257  $PA_s$  are the concomitant accuracies in stratum  $s$ , and  $W_s$  is the sampling weight of strata.

## 258 **5 Results**

### 259 **5.1 Water bodies in tundra and boreal forests of North America**

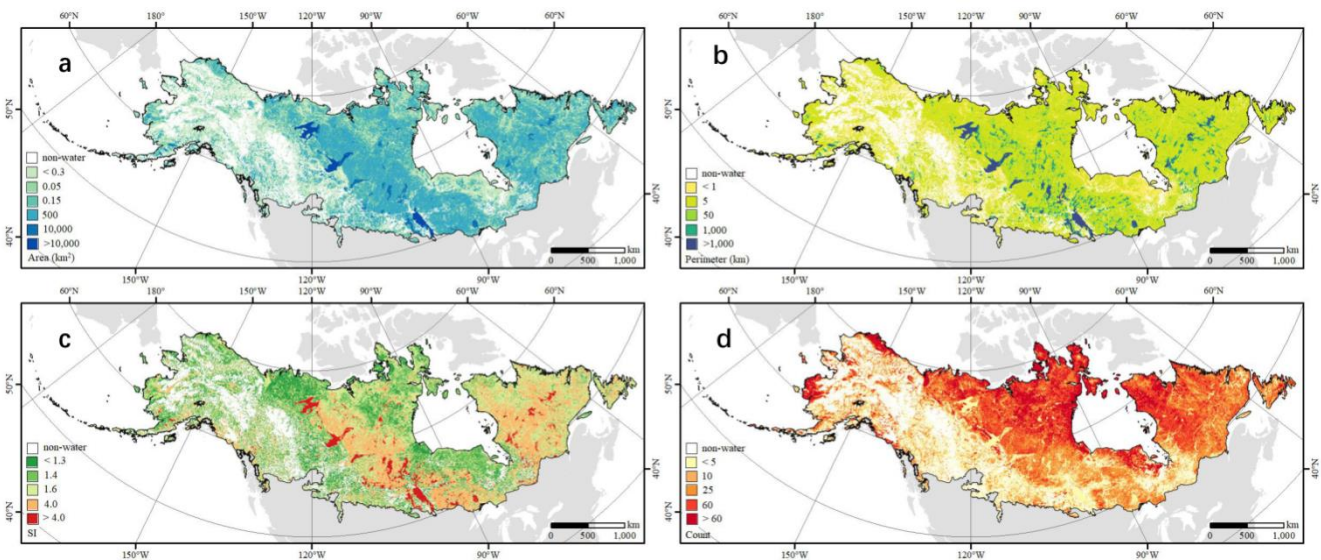
260 More than 6.47 million (6,474,051) surface water bodies were identified in the tundra and boreal forests of North America,  
261 while 90.3% of these water bodies (5,844,921) were smaller than 0.1 km<sup>2</sup>. Those water bodies covered more than 0.8 million  
262 km<sup>2</sup>, ~10.3% of the study area (Figure 6). The average size and perimeter of the identified water bodies were 0.12 km<sup>2</sup> and  
263 1.01 km, respectively, and their average  $SI$  was 1.41.



264

265 **Figure 6: Percent of surface water (5 km × 5km grid) produced by aggregating the water extent for the tundra and boreal forests of**  
 266 **North America as calculated using the WBD-NAHL dataset.**

267 All of the morphological indicators, including area, perimeter, and *SI*, of the identified water bodies showed great heterogeneity  
 268 across the region (Figure 7). In general, the tundra biome consists of a large number of densely packed small water bodies with  
 269 regular shapes. In contrast, the boreal forest biome consists of a large number of large water bodies with complex shapes. The  
 270 number of identified water bodies in the tundra (3.24 million) and boreal forests (3.23 million) were nearly identical. However,  
 271 the water extent in the boreal forest (0.57 million km<sup>2</sup>; 71% of total water area) is more than twice that found in the tundra  
 272 (0.23 million km<sup>2</sup>; 29% of the total water area), indicating that the average size of water bodies in the boreal area are larger  
 273 than those in the tundra. This finding was confirmed by reviewing the water body perimeters for the two biomes. The average  
 274 perimeter of water bodies in boreal forests was 1.2 km, compared to a much smaller 0.8 km average perimeter for water bodies  
 275 in the tundra. The average *SI* for water bodies in the boreal was 1.45, longer than the 1.37 average *SI* for the tundra water  
 276 bodies, suggesting that the boreal water bodies generally have much more complex shorelines, while the tundra water bodies  
 277 are more circular.



278

279 **Figure 7: The aggregated distribution of area (a), perimeter (b), and *SI* (c), and the number (d) of the identified water bodies in the**  
 280 **study area. The values at each 5 km × 5 km pixel in the grid were calculated by selecting the intersecting water bodies and then**  
 281 **either counting or calculating the mean of the targeted parameter (e.g., area, *SI*, and perimeter) of these selected water bodies.**

282

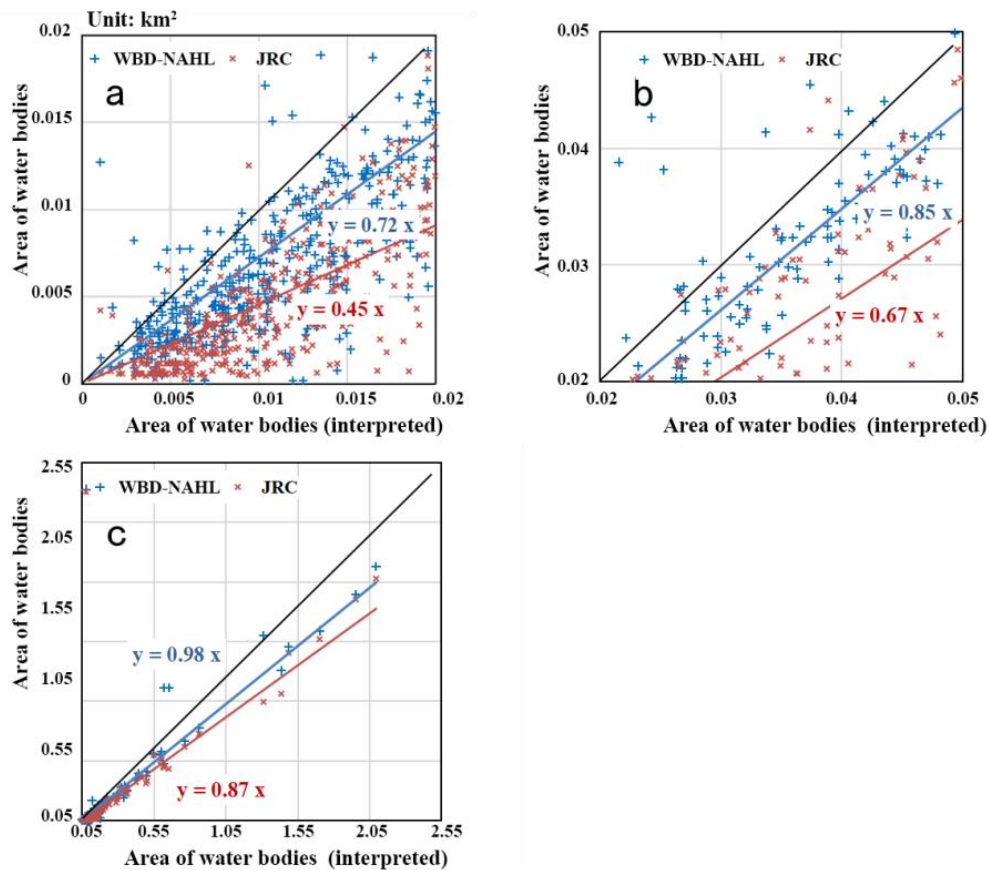
283 Inland water in the region is mainly concentrated in the Canadian Shield, i.e., about 0.73 million km<sup>2</sup> of water (92% of water  
 284 extent in the study region). In addition, most large water bodies were located in the Canadian Shield, including 90% of the

285 identified large water bodies (sizes  $\geq 1\text{km}^2$ ). The shorelines of the water bodies in the Canadian Shield were also more  
286 complex than those in other areas, especially south of the Laurentian Plateau near the Great Lakes.

## 287 5.2 Accuracy assessment

288 The overall accuracy of the WBD-NAHL's water extent was 96.36%, while the producer's accuracy was 99.9%, and the user's  
289 accuracy was 96.36%. Misclassifications were primarily found in shadows of the Mackenzie Mountains, where the east-west  
290 high-elevation mountain range cast constant shadows on the northern slopes.

291 Both the JRC and WBD-NAHL datasets accurately identified the size of larger water bodies. For mixed water pixels, the area  
292 estimates of both datasets were more conservative than the reference data. However, the WBD-NAHL dataset performed better  
293 than the JRC. The advantage of the WBD-NAHL was demonstrated for smaller water bodies (Figure 8). For small water bodies  
294 (size  $\leq 0.02\text{ km}^2$ ), the average area of the WBD-NAHL water bodies was 72% of those manually digitized over high-  
295 resolution Google Earth images, compared to only 45% with the water area detected by the JRC (Figure 8a). For medium  
296 water bodies (between  $0.02\text{ km}^2$  and  $0.05\text{ km}^2$ ), the average area of WBD-NAHL water bodies was about 85% times that of  
297 manually digitized water bodies, compared to 67% with the water area detected by the JRC (Figure 8b). For water bodies  
298 larger than  $0.05\text{ km}^2$ , the water areas of WBD-NAHL were highly consistent (98%) with that of manually digitized, while the  
299 water area of JRC was slightly lower (about 87%) for water bodies in the category (Figure 8c).

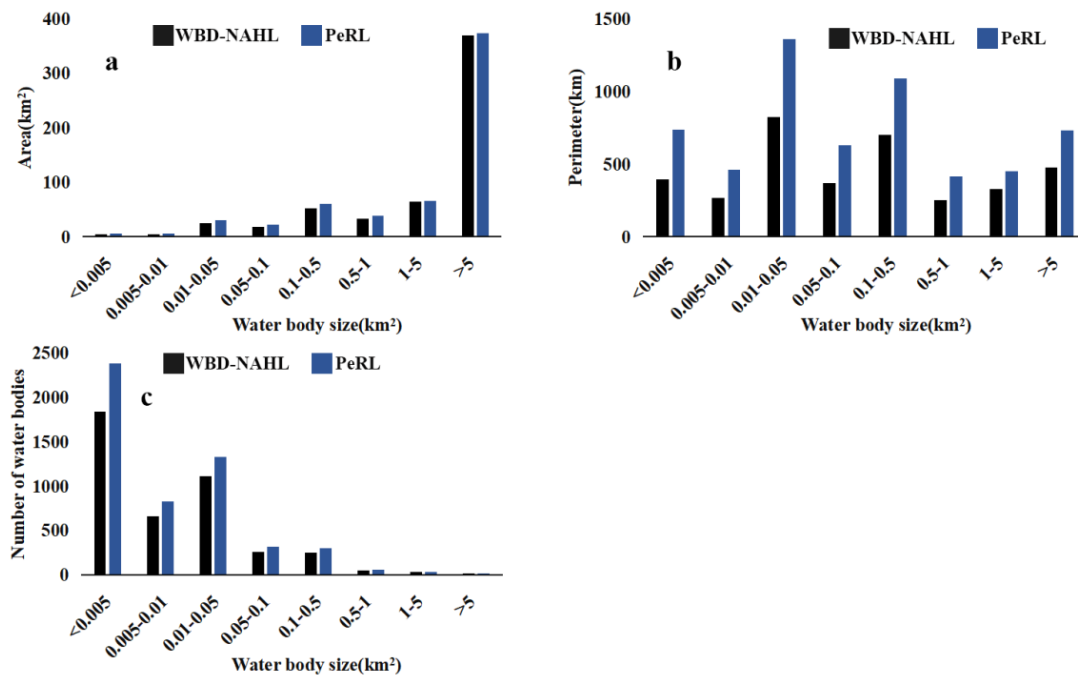


300

301 **Figure 8:** Comparisons of the water body area identified by the JRC, WBD-NAHL, and interpreted water maps. The 1:1 lines are  
302 in black. The red crosses represent the JRC water bodies, and the blue pluses represent the WBD-NAHL water bodies in comparison  
303 with the manually interpreted water bodies. The water bodies are compared in groups of sizes, i.e., (a) small water bodies with sizes  
304  $< 0.02\text{ km}^2$ ; (b) medium water bodies with sizes between  $0.02\text{ km}^2$  and  $0.05\text{ km}^2$ ; (c) large water bodies with sizes  $> 0.05\text{ km}^2$ . The  $R^2$   
305 for the WBD-NAHL and JRC identified water bodies were similar, i.e., 0.6 for small water bodies, 0.5 for medium water bodies, and  
306 0.9 for large water bodies.

307

308 The comparison between the water bodies identified by WBD-NAHL and PeRL were largely consistent for the derived  
 309 indicators of water area, perimeter, and number (Figure 9). Linear correlations between the water bodies identified by WBD-  
 310 NAHL and PeRL had  $R^2$  higher than 0.99 for all three indicators. The slopes of the linear regressions indicated that the water  
 311 area showed the least bias when compared to PeRL (slope=0.98), followed by the number of water bodies (slope=0.78), and  
 312 finally the perimeter of the water bodies (slope=0.62).



313  
 314 **Figure 9: Area, perimeter, and number of water bodies identified by the PeRL and WBD-NAHL datasets.**

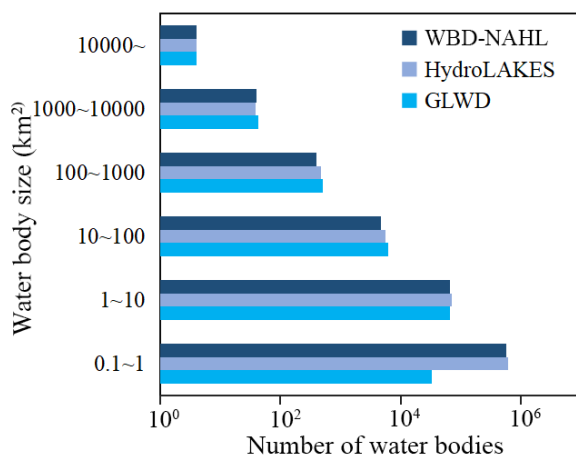
## 315 6 Discussion

### 316 6.1 A high-resolution water body dataset for the continental tundra and boreal

317 The WBD-NAHL dataset provides the first known delineation of water bodies at 10-m resolution for the continental tundra  
 318 and boreal forest of North America, which is one of the highest concentrations of the global inland water, especially the small  
 319 sized water bodies. The dataset not only maps the extent of inland water during 2019 but also identifies the water bodies and  
 320 their morphological metrics, which are critical for understanding and modeling freshwater lentic ecosystems (Downing, 2009;  
 321 Heathcote et al., 2015; Kuhn and Butman, 2021; MacIntyre et al., 2009; Muster et al., 2013). The WBD-NAHL was produced  
 322 using Sentinel-2 satellite data to take advantage of the high resolution and 2-3-day revisit time of Sentinel-2 satellites. Sentinel-  
 323 2's revisit time allows the WBD-NAHL to have sufficient observations during the snow-free season, which is critical for  
 324 mapping inland surface water in this high latitude region with long periods of snow coverage.

325 The WBD-NAHL's 10-m resolution enabled detecting water bodies as small as 0.001 km<sup>2</sup>. The validation showed that the  
 326 WBD-NAHL dataset had high overall accuracy and significantly improved upon the ability of the existing global JRC water  
 327 maps for detecting small water (e.g., smaller than 0.006 km<sup>2</sup>) than the existing global JRC water maps. These small water  
 328 bodies consist of nearly half the total water bodies in the tundra and boreal forest regions of North America, and generally  
 329 experience faster cycling of water, material, and energy than larger water bodies (Winslow et al., 2014; Carroll et al., 2011;  
 330 Messenger et al., 2016). The improved WBD-NAHL dataset may provide more accurate inputs for hydrological estimates,  
 331 which are vital components for understanding and modeling the pan-Arctic hydrological, biochemical, and energy cycling.

332 The higher resolution of WBD-NAHL also provides the ability to delineate the number, area, and shoreline complexity of  
 333 water bodies. Our comparison confirmed that WBD-NAHL-derived water areas and shorelines were similar to those from the  
 334 regional 5-m or finer resolution PeRL dataset. Meanwhile, the number of water bodies identified in WBD-NAHL was  
 335 consistent with those of other datasets, including HydroLAKES and GLWD (Figure 10). The number of water bodies larger  
 336 than 1 km<sup>2</sup> was roughly identical for WBD-NAHL, HydroLAKES, and GLWD. For water bodies between 0.1 and 1 km<sup>2</sup>,  
 337 WBD-NAHL and HydroLAKES reported similar numbers (Figure 10), but the number reported by GLWD was considerably  
 338 lower, suggesting that the omission error of GLWD was higher for water bodies smaller than 1 km<sup>2</sup>, as noted by Lehner and  
 339 Döll (2004). Unfortunately, both the HydroLAKES and GLWD datasets only provide records for water bodies larger than 0.1  
 340 km<sup>2</sup> (Messenger et al., 2016; Lehner and Döll, 2004), and are thus missing records for what we estimate to be 90% of the total  
 341 number of water bodies in the region. The WBD-NAHL is able to extend these indicators to much smaller water bodies than  
 342 HydroLAKES and GLWD, providing a much more complete record of water bodies in the region. This estimate of the number  
 343 and extent of small water bodies can improve our understanding of continental freshwater sources, stressing the importance of  
 344 small water bodies in continental biochemical and energy cycling, potentially correcting a misconception that large lakes are  
 345 most important (Downing, 2010).

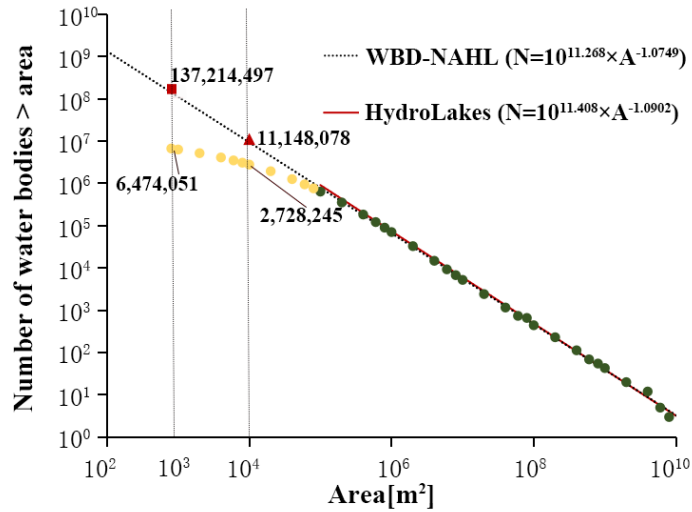


346

347 **Figure 10: Comparing the number of water bodies identified by WBD-NAHL and other datasets based on size.**

348 **6.2 Distribution of the water bodies**

349 An empirical power-law distribution was found between lake areas and lake numbers (Messenger et al., 2016; Downing et al.,  
 350 2006), and the distribution was applied to estimate the number of small lakes, which were used for estimating greenhouse gas  
 351 emissions (Holgerson et al., 2016). According to the power-law distribution and HydroLAKES, the number of water bodies  
 352 larger than 0.1 km<sup>2</sup> was estimated to be about 798,895, which was close to the 629,130 water bodies reported by WBD-NAHL  
 353 (Figure 11). However, the number of water bodies sized between 0.1 and 0.01 km<sup>2</sup> was estimated to be about 10.2 million, 4.8  
 354 times higher than estimated by WBD-NAHL. Furthermore, the water bodies sized between 0.01-0.001 km<sup>2</sup> were estimated to  
 355 be about 126.1 million, 33.6 times higher than what was estimated by WBD-NAHL, suggesting that the power-law distribution  
 356 significantly overestimates the number of small lakes. A similar finding was reported by Seekell et al. (2016). Estimating the  
 357 number small water bodies using a power-law distribution could introduce considerable uncertainties in the estimation of the  
 358 contribution of small water bodies to greenhouse gas emissions. Accurately identifying small water bodies could correct this  
 359 overestimation and improve greenhouse gas emission estimates (Holgerson et al., 2016).



360

361 **Figure 11: Distribution of the total numbers of water bodies in relation to the area of North American tundra and boreal forests**  
 362 **water bodies. The circles represent the number of water bodies provided by WBD-NAHL. The black line is the power-law**  
 363 **distribution modeled using water bodies > 0.1 km<sup>2</sup> from WBD-NAHL. The red line is the power-law distribution modeled using**  
 364 **HydroLAKES in the study region. The red triangle and square represent, respectively, the extrapolated number of water bodies >**  
 365 **0.01 km<sup>2</sup> and > 0.001 km<sup>2</sup> based on the power-law distribution modeled from HydroLAKES.**

366

367 The largest and most complex water bodies are distributed primarily in the Canadian Shield. These lakes in the Canadian  
 368 Shield formed through processes such as erosion and glaciation (Smith et al., 2007). Erosion and glaciation formed water  
 369 bodies with complex shapes, which may contribute to the higher *SI* (1.48) reported by the WBD-NAHL for the region. During  
 370 the most recent Wisconsin glaciation, the Canadian Shield was covered by the Laurentide Ice Sheet, a giant, 3-km thick expanse  
 371 of ice. When the ice sheet retreated north, it carved out the five Great Lakes as well as thousands of small lakes throughout the  
 372 Canadian Shield (Dyke and Prest, 1987). Currently, 92% of the water extent in the tundra and boreal forests are distributed in  
 373 this particular region. For example, the largest lake in the region - Great Bear Lake - has a surface area of 30,227 km<sup>2</sup> with a  
 374 long, complex shoreline (the perimeter is 5,705 km and the *SI* of the lake is 9.3). It was formed by ice erosion during the  
 375 Pleistocene (Johnson, 1975).

376 The tundra, on the other hand, has a large number of small, regularly-shaped water bodies, which could be related to the thick  
 377 peatland and thermokarst landscape. Over the past few decades, numerous thermokarst lakes have been experiencing dramatic  
 378 changes, which are considered an indicator of permafrost degradation (Smith et al., 2005; Karlsson et al., 2012, 2014). The  
 379 small thermokarst lakes were also found to experience stronger changes than larger lakes (Karlsson et al., 2014; Carroll and  
 380 Loboda, 2017). Monitoring water extent without discriminating by lake size does not accurately reflect these changes in small  
 381 lakes due to the area dominance of large lakes. Additionally, the small thermokarst lakes are the primary source of permafrost  
 382 carbon emissions (Kuhn et al., 2018; Walter Anthony et al., 2016; Yvon-Durocher et al., 2017) and small water bodies were  
 383 found to be a major source of uncertainty greenhouse gas emission estimates (Holgerson and Raymond, 2016). The WBD-  
 384 NAHL dataset could provide critical information for investigating thermokarst lakes, especially small thermokarst lakes and  
 385 ponds, and estimating their effects on carbon emission and permafrost sustainability in the tundra and boreal forests of North  
 386 America. As reported by the analysis of WBD-NAHL, 3.24 million small water bodies were found in the tundra in 2019, with  
 387 an average size of 0.07 km<sup>2</sup> and average *SI* of 1.37, much smaller than the *SI* of boreal lakes. Teshekpuk Lake is the largest  
 388 thermokarst lake in the world and a relatively smooth shoreline (*SI* = 5.4), considerably smaller than the *SI* of Great Bear Lake  
 389 in the boreal region (Markon and Derksen, 1994).

390 The biome-based analysis provided insights into the distribution of water body shapes across the study area; however, more  
391 complex relationships can be found between the shapes and the surface geology of the water bodies. For example, circular-  
392 shaped lakes can be found in regions with thick overburden – possibly as a result from remaining unglaciated, from aeolian  
393 deposits or from rising from the sea bottom through isostatic rebound; These circular-shaped lakes can be found in regions  
394 with thick moraines or widespread peatlands in the boreal Hudson Bay lowlands and the Mackenzie River Basin. The high-  
395 resolution WBD-NAHL could help further explore the distribution of water bodies by size and shape.

### 396 **6.3 Limitations**

397 The data and methods used to derive the 10-m resolution WBD-NAHL dataset are able to detect water bodies smaller than the  
398 30-m or coarser-resolution satellite-derived datasets, but have difficulty identifying water bodies smaller than 0.001 km<sup>2</sup>. This  
399 limitation can be further improved by incorporating higher resolution satellite data, such as from Planet, WorldView,  
400 QuickBird, and Gaofen (Veremeeva and Günther, 2017; Sun et al., 2020; Watson et al., 2016; Andresen and Lougheed, 2015).  
401 Limit errors in the satellite data provide substantial sources of uncertainty, including an inability to separate rivers and streams  
402 because the resolution is too coarse, bias in estimates of water extent resulting from temporal gaps in the data, and  
403 misclassifications resulting from spectral resolution. The misclassifications impacted by terrain (e.g., mountain shadows) still  
404 exist even though they have been substantially reduced during data processing. Further processing may be possible to further  
405 reduce these errors.

406 The WBD-NAHL dataset was produced based on Sentinel-2 data acquired in the summer of 2019 and represents the  
407 distribution of surface water in the corresponding year. The mean total precipitation in 2019 in the region was 438.5 mm,  
408 which was close to the historical average from 2010 to 2019 (mean: 435.9 mm, standard deviation: 11.5 mm) (Huffman et al.  
409 2019). Although 2019 can be considered a normal year of the past decade in terms of precipitation, the spatial extent of high-  
410 latitude water bodies, especially smaller water bodies, can still vary significantly both inter- and intra-annually locally.  
411 Nevertheless, it would be interesting to explore water bodies' changes using observations from multiple years. Further efforts  
412 can be carried out to produce an inland water dataset for multiple time periods using these methods to capture the seasonal and  
413 multi-year dynamics of inland water in the region. The WBD-NAHL dataset focused on the tundra and boreal forest regions  
414 of North America. The methodology can be extended to Eurasia to provide a complete representation of the biomes.

### 415 **7 Data availability**

416 This WBD-NAHL dataset can be accessed via the website of the National Tibetan Plateau/Third Pole Environment Data Center  
417 (TPDC, <http://data.tpdc.ac.cn>): DOI: 10.11888/Hydro.tpdc.271021 (Feng et al., 2020). The dataset is provided in ESRI  
418 Geodatabase format. The volume of this dataset is about 1.5 GB.

### 419 **8 Conclusions**

420 This study presents an inland surface water body dataset for the North American high latitudes. The WBD-NAHL dataset was  
421 generated using Sentinel-2 data with machine learning methods and an object-based algorithm. Three morphological metrics  
422 (area, perimeter, and *SI*) were calculated for each water body. The accuracy of the dataset was carefully assessed with respect  
423 to detecting inland surface water extent (or pixel level) and identifying water bodies. The dataset's overall accuracy for water  
424 extent reached 96.36%. In addition, the WBD-NAHL showed a high consistency with high-resolution images in terms of water  
425 area, perimeter, and quantity.

426 To our knowledge, the WBD-NAHL dataset provided the most complete inventory of inland surface water bodies for the  
427 tundra and boreal forest regions of North America. Overall, 6.47 million water bodies were identified, covering 10.3% of the  
428 region. Small water bodies dominate the region, as ~90.3% have an area smaller than 0.1 km<sup>2</sup>. The WBD-NAHL indicates that  
429 the tundra biome is dominated by densely distributed small water bodies with regular shapes (the average *SI* was 1.37), while  
430 the boreal forest biome is dominated by large water bodies with complex shapes (the average *SI* was 1.45). The WBD-NAHL  
431 is expected to be able to provide supporting data for modeling hydrologic, biochemical, and energy cycling in these areas.

## 432 Acknowledgements

433 This work was supported by Basic Science Center for Tibetan Plateau Earth System (BSCTPES, NSFC project No. 41988101)

## 434 Reference

- 435 Andresen, C. G. and Lougheed, V. L.: Disappearing Arctic tundra ponds: Fine-scale analysis of surface hydrology in drained  
436 thaw lake basins over a 65 year period (1948–2013), *J. Geophys. Res. Biogeosciences*, 120, 466–479, 2015.
- 437 Biskaborn, B. K., Smith, S. L., Noetzli, J., Matthes, H., Vieira, G., Streletskiy, D. A., Schoeneich, P., Romanovsky, V. E.,  
438 Lewkowicz, A. G., and Abramov, A.: Permafrost is warming at a global scale, *Nat. Commun.*, 10, 1–11, 2019.
- 439 Breiman, L.: Random forests, *Mach. Learn.*, 45, 5–32, 2001.
- 440 Carlson, T. N. and Ripley, D. A.: On the relation between NDVI, fractional vegetation cover, and leaf area index, *Remote  
441 Sens. Environ.*, 62, 241–252, 1997.
- 442 Carpenter, S. R.: Lake geometry: implications for production and sediment accretion rates, *J. Theor. Biol.*, 105, 273–286, 1983.
- 443 Carroll, M. and Loboda, T.: Multi-Decadal Surface Water Dynamics in North American Tundra, *Remote Sens.*, 9,  
444 <https://doi.org/10.3390/rs9050497>, 2017.
- 445 Carroll, M. L., Townshend, J. R. G., DiMiceli, C. M., Loboda, T., and Sohlberg, R. A.: Shrinking lakes of the Arctic: Spatial  
446 relationships and trajectory of change, *Geophys. Res. Lett.*, 38, 2011.
- 447 Cooley, S. W., Smith, L. C., Ryan, J. C., Pitcher, L. H., and Pavelsky, T. M.: Arctic-Boreal Lake Dynamics Revealed Using  
448 CubeSat Imagery, *Geophys. Res. Lett.*, 46, 2111–2120, <https://doi.org/10.1029/2018gl081584>, 2019.
- 449 Danielson, J. J. and Gesch, D. B.: Global multi-resolution terrain elevation data 2010 (GMTED2010), 2011.
- 450 Downing, J. A.: Global limnology: Up-scaling aquatic services and processes to planet Earth, *Int. Ver. Für Theor. Angew.  
451 Limnol. Verhandlungen*, 30, 1149–1166, 2009.
- 452 Downing, J. A.: Emerging global role of small lakes and ponds: little things mean a lot, *Limnetica*, 29, 0009–0024, 2010.
- 453 Dranga, S. A., Hayles, S., and Gajewski, K.: Synthesis of limnological data from lakes and ponds across Arctic and Boreal  
454 Canada, *Arct. Sci.*, 4, 167–185, <https://doi.org/10.1139/as-2017-0039>, 2017.
- 455 Du, J., Kimball, J. S., Jones, L. A., and Watts, J. D.: Implementation of satellite based fractional water cover indices in the  
456 pan-Arctic region using AMSR-E and MODIS, *Remote Sens. Environ.*, 184, 469–481,  
457 <https://doi.org/10.1016/j.rse.2016.07.029>, 2016.
- 458 Dyke, A. and Prest, V.: Late Wisconsinan and Holocene history of  
459 the Laurentide ice sheet, *Géographie Phys. Quat.*, 41, 237–263, 1987.
- 459 Fayne, J. V., Smith, L. C., Pitcher, L. H., Kyzivat, E. D., Cooley, S. W., Cooper, M. G., Denbina, M. W., Chen, A. C., Chen,  
460 C. W., and Pavelsky, T. M.: Airborne observations of arctic-boreal water surface elevations from AirSWOT Ka-Band  
461 InSAR and LVIS LiDAR, *Environ. Res. Lett.*, 15, 105005, 2020.
- 462 Feng, M., Sexton, J. O., Channan, S., and Townshend, J. R.: A global, high-resolution (30-m) inland water body dataset for  
463 2000: first results of a topographic–spectral classification algorithm, *Int. J. Digit. Earth*, 9, 113–133,  
464 <https://doi.org/10.1080/17538947.2015.1026420>, 2015.
- 465 Feng, M., Sui, Y.: High resolution inland surface water dataset for the tundra and boreal in North America,  
466 <https://doi.org/10.11888/Hydro.tpcd.271021>, 24 October 2020.
- 467 Forkel, M., Carvalhais, N., Rödenbeck, C., Keeling, R., Heimann, M., Thonicke, K., Zaehle, S., and Reichstein, M.: Enhanced  
468 seasonal CO<sub>2</sub> exchange caused by amplified plant productivity in northern ecosystems, *Science*, 351, 696–699, 2016.
- 469 Glińska-Lewczuk, K.: Water quality dynamics of oxbow lakes in young glacial landscape of NE Poland in relation to their  
470 hydrological connectivity, *Ecol. Eng.*, 35, 25–37, <https://doi.org/10.1016/j.ecoleng.2008.08.012>, 2009.
- 471 Graversen, R. G., Mauritsen, T., Tjernström, M., Källén, E., and Svensson, G.: Vertical structure of recent Arctic warming,  
472 *Nature*, 451, 53–56, 2008.
- 473 Grosse, G., Jones, B., and Arp, C.: 8.21 Thermokarst Lakes, Drainage, and Drained Basins, in: *Treatise on Geomorphology*,  
474 edited by: Shroder, J. F., Academic Press, San Diego, 325–353, <https://doi.org/10.1016/B978-0-12-374739-6.00216-5>, 2013.
- 475 Han-Qiu, X.: A study on information extraction of water body with the modified normalized difference water index (MNDWI),  
476 *J. Remote Sens.*, 5, 589–595, 2005.
- 477 Heathcote, A. J., del Giorgio, P. A., and Prairie, Y. T.: Predicting bathymetric features of lakes from the topography of their  
478 surrounding landscape, *Can. J. Fish. Aquat. Sci.*, 72, 643–650, 2015.
- 479



480 Higgins, S., Desjardins, C., Drouin, H., Hrenchuk, L., and Van der Sanden, J.: The role of climate and lake size in regulating  
481 the ice phenology of boreal lakes, *J. Geophys. Res. Biogeosciences*, 126, e2020JG005898, 2021.

482 Holgerson, M. A. and Raymond, P. A.: Large contribution to inland water CO<sub>2</sub> and CH<sub>4</sub> emissions from very small ponds,  
483 *Nat. Geosci.*, 9, 222–226, <https://doi.org/10.1038/ngeo2654>, 2016.

484 Huffman, G.J., E.F. Stocker, D.T. Bolvin, E.J. Nelkin, Jackson Tan (2019), GPM IMERG Final Precipitation L3 1 month 0.1  
485 degree x 0.1 degree V06, Greenbelt, MD, Goddard Earth Sciences Data and Information Services Center (GES DISC),  
486 Accessed: 2022.5.24, 10.5067/GPM/IMERG/3B-MONTH/06

487 van Huissteden, J., Berrittella, C., Parmentier, F. J. W., Mi, Y., Maximov, T. C., and Dolman, A. J.: Methane emissions from  
488 permafrost thaw lakes limited by lake drainage, *Nat. Clim. Change*, 1, 119–123, <https://doi.org/10.1038/nclimate1101>,  
489 2011.

490 Isikdogan, F., Bovik, A. C., and Passalacqua, P.: Surface Water Mapping by Deep Learning, *IEEE J. Sel. Top. Appl. Earth*  
491 *Obs. Remote Sens.*, 10, 4909–4918, <https://doi.org/10.1109/JSTARS.2017.2735443>, 2017.

492 Jiang, X., Zheng, P., Cao, L., and Pan, B.: Effects of long-term floodplain disconnection on multiple facets of lake fish  
493 biodiversity: Decline of alpha diversity leads to a regional differentiation through time, *Sci. Total Environ.*, 763,  
494 144177, 2021.

495 Johannessen, O. M., Bengtsson, L., Miles, M. W., Kuzmina, S. I., Semenov, V. A., Alekseev, G. V., Nagurnyi, A. P., Zakharov,  
496 V. F., Bobylev, L. P., and Pettersson, L. H.: Arctic climate change: observed and modelled temperature and sea-ice  
497 variability, *Tellus Dyn. Meteorol. Oceanogr.*, 56, 328–341, 2004.

498 Johnson, L.: The Great Bear Lake: its place in history, *Arctic*, 28, 231–244, 1975.

499 Karlsson, J., Lyon, S., and Destouni, G.: Temporal Behavior of Lake Size-Distribution in a Thawing Permafrost Landscape in  
500 Northwestern Siberia, *Remote Sens.*, 6, 621–636, <https://doi.org/10.3390/rs6010621>, 2014.

501 Karlsson, J. M., Lyon, S. W., and Destouni, G.: Thermokarst lake, hydrological flow and water balance indicators of permafrost  
502 change in Western Siberia, *J. Hydrol.*, 464–465, 459–466, <https://doi.org/10.1016/j.jhydrol.2012.07.037>, 2012.

503 King, K. B., Bremigan, M. T., Infante, D., and Cheruvelil, K. S.: Surface water connectivity affects lake and stream fish species  
504 richness and composition, *Can. J. Fish. Aquat. Sci.*, 78, 433–443, 2021.

505 Kuhn, C. and Butman, D.: Declining greenness in Arctic-boreal lakes, *Proc. Natl. Acad. Sci.*, 118, 2021.

506 Kuhn, M., Lundin, E. J., Giesler, R., Johansson, M., and Karlsson, J.: Emissions from thaw ponds largely offset the carbon  
507 sink of northern permafrost wetlands, *Sci. Rep.*, 8, 9535, <https://doi.org/10.1038/s41598-018-27770-x>, 2018.

508 Laird, N. F., Walsh, J. E., and Kristovich, D. A.: Model simulations examining the relationship of lake-effect morphology to  
509 lake shape, wind direction, and wind speed, *Mon. Weather Rev.*, 131, 2102–2111, 2003.

510 Langer, M., Westermann, S., Boike, J., Kirillin, G., Grosse, G., Peng, S., and Krinner, G.: Rapid degradation of permafrost  
511 underneath waterbodies in tundra landscapes—Toward a representation of thermokarst in land surface models, *J.*  
512 *Geophys. Res. Earth Surf.*, 121, 2446–2470, <https://doi.org/10.1002/2016jfr003956>, 2016.

513 Laske, S. M., Rosenberger, A. E., Wipfli, M. S., and Zimmerman, C. E.: Surface water connectivity controls fish food web  
514 structure and complexity across local- and meta-food webs in Arctic Coastal Plain lakes, *Food Webs*, 21,  
515 <https://doi.org/10.1016/j.fooweb.2019.e00123>, 2019.

516 Lehner, B. and Döll, P.: Development and validation of a global database of lakes, reservoirs and wetlands, *J. Hydrol.*, 296,  
517 1–22, <https://doi.org/10.1016/j.jhydrol.2004.03.028>, 2004.

518 Li, X., Che, T., Li, X., Wang, L., Duan, A., Shangguan, D., Pan, X., Fang, M., and Bao, Q.: CASEarth poles: big data for the  
519 three poles, *Bull. Am. Meteorol. Soc.*, 101, E1475–E1491, 2020.

520 Lindgren, P. R., Farquharson, L. M., Romanovsky, V. E., and Grosse, G.: Landsat-based lake distribution and changes in  
521 western Alaska permafrost regions between the 1970s and 2010s, *Environ. Res. Lett.*, 16, 025006,  
522 <https://doi.org/10.1088/1748-9326/abd270>, 2021.

523 MacIntyre, S., Fram, J. P., Kushner, P. J., Bettez, N. D., O'Brien, W., Hobbie, J., and Kling, G. W.: Climate-related variations  
524 in mixing dynamics in an Alaskan arctic lake, *Limnol. Oceanogr.*, 54, 2401–2417, 2009.

525 Markon, C. J. and Derksen, D. V.: Identification of tundra land cover near Teshekpuk Lake, Alaska using SPOT satellite data,  
526 *Arctic*, 222–231, 1994.

527 McCullough, I. M., King, K. B. S., Stachelek, J., Diaz, J., Soranno, P. A., and Cheruvelil, K. S.: Applying the patch-matrix  
528 model to lakes: a connectivity-based conservation framework, *Landsc. Ecol.*, 34, 2703–2718,  
529 <https://doi.org/10.1007/s10980-019-00915-7>, 2019.

530 McFeeters, S. K.: The use of the Normalized Difference Water Index (NDWI) in the delineation of open water features, *Int. J.*  
531 *Remote Sens.*, 17, 1425–1432, 1996.

532 Messenger, M. L., Lehner, B., Grill, G., Nedeva, I., and Schmitt, O.: Estimating the volume and age of water stored in global  
533 lakes using a geo-statistical approach, *Nat. Commun.*, 7, 1–11, 2016.

534 Meyer, M. F., Labou, S. G., Cramer, A. N., Brousil, M. R., and Luff, B. T.: The global lake area, climate, and population  
535 dataset, *Sci Data*, 7, 174, <https://doi.org/10.1038/s41597-020-0517-4>, 2020.

536 Muster, S., Heim, B., Abnizova, A., and Boike, J.: Water body distributions across scales: A remote sensing based comparison  
537 of three arctic tundra wetlands, *Remote Sens.*, 5, 1498–1523, 2013.

538 Napiórkowski, Bąkowska, Mrozińska, Szymańska, Kolarova, and Obolewski: The Effect of Hydrological Connectivity on the  
539 Zooplankton Structure in Floodplain Lakes of a Regulated Large River (the Lower Vistula, Poland), *Water*, 11,  
540 <https://doi.org/10.3390/w11091924>, 2019.

541 Nitze, I., Cooley, S. W., Duguay, C. R., Jones, B. M., and Grosse, G.: The catastrophic thermokarst lake drainage events of  
542 2018 in northwestern Alaska: Fast-forward into the future, *The Cryosphere*, 14, 4279–4297, 2020.

543 Olefeldt, D., Goswami, S., Grosse, G., Hayes, D., Hugelius, G., Kuhry, P., McGuire, A. D., Romanovsky, V. E., Sannel, A. B.  
544 K., Schuur, E. a. G., and Turetsky, M. R.: Circumpolar distribution and carbon storage of thermokarst landscapes,  
545 *Nat. Commun.*, 7, 13043, <https://doi.org/10.1038/ncomms13043>, 2016

546 Pachauri, R. K. and Reisinger, A.: IPCC fourth assessment report, IPCC Geneva, 2007, 2007.

547 Pekel, J.-F., Cottam, A., Gorelick, N., and Belward, A. S.: High-resolution mapping of global surface water and its long-term  
548 changes, *Nature*, 540, 418–422, <https://doi.org/10.1038/nature20584>, 2016.

549 Pickens, A. H., Hansen, M. C., Hancher, M., Stehman, S. V., Tyukavina, A., Potapov, P., Marroquin, B., and Sherani, Z.:  
550 Mapping and sampling to characterize global inland water dynamics from 1999 to 2018 with full Landsat time-series,  
551 *Remote Sens. Environ.*, 243, <https://doi.org/10.1016/j.rse.2020.111792>, 2020.

552 Ritter, M. E.: The physical environment: An introduction to physical geography, Date Visit. July, 25, 2008, 2006.

553 Sandlund, O. T., Eloranta, A. P., Borgström, R., Hesthagen, T., Johnsen, S. I., Museth, J., and Rognerud, S.: The trophic niche  
554 of Arctic charr in large southern Scandinavian lakes is determined by fish community and lake morphometry,  
555 *Hydrobiologia*, 783, 117–130, <https://doi.org/10.1007/s10750-016-2646-5>, 2016.

556 Schilder, J., Bastviken, D., van Hardenbroek, M., Kankaala, P., Rinta, P., Stötter, T., and Heiri, O.: Spatial heterogeneity and  
557 lake morphology affect diffusive greenhouse gas emission estimates of lakes, *Geophys. Res. Lett.*, 40, 5752–5756,  
558 2013.

559 Serikova, S., Pokrovsky, O. S., Laudon, H., Krickov, I., Lim, A. G., Manasypov, R. M., and Karlsson, J.: High carbon  
560 emissions from thermokarst lakes of Western Siberia, *Nat. Commun.*, 10, 1–7, 2019.

561 Serreze, M. C. and Francis, J. A.: The Arctic amplification debate, *Clim. Change*, 76, 241–264, 2006.

562 Sharma, S., Blagrove, K., Magnuson, J. J., O'Reilly, C. M., Oliver, S., Batt, R. D., Magee, M. R., Straile, D., Weyhenmeyer,  
563 G. A., Winslow, L., and Woolway, R. I.: Widespread loss of lake ice around the Northern Hemisphere in a warming  
564 world, *Nat. Clim. Change*, 9, 227–231, <https://doi.org/10.1038/s41558-018-0393-5>, 2019.

565 Smith, L. C., Sheng, Y., MacDonald, G. M., and Hinzman, L. D.: Disappearing arctic lakes, *Science*, 308, 1429–1429, 2005.

566 Smith, L. C., Sheng, Y., and MacDonald, G. M.: A first pan-Arctic assessment of the influence of glaciation, permafrost,  
567 topography and peatlands on northern hemisphere lake distribution, *Permafr. Periglac. Process.*, 18, 201–208,  
568 <https://doi.org/10.1002/ppp.581>, 2007.

569 Sun, J., Wang, G., He, G., Pu, D., Jiang, W., Li, T., and Niu, X.: Study on the Water Body Extraction Using GF-1 Data Based  
570 on Adaboost Integrated Learning Algorithm, *Int. Arch. Photogramm. Remote Sens. Spat. Inf. Sci.*, 42, 641–648, 2020.

571 Van Gerven, M. and Bohte, S.: Artificial neural networks as models of neural information processing, *Front. Comput.*  
572 *Neurosci.*, 11, 114, 2017.

573 Veremeeva, A. A. and Günther, F.: Thermokarst lake and baydzherakh area changes on Yedoma uplands, Yakutian coastal  
574 lowlands: repeat inventory using high resolution imagery, 2017.

575 Vaideliene, A., & Michailov, N. (2008). Dam influence on the river self-purification. 748–757.

576 Walter Anthony, K., Daanen, R., Anthony, P., Schneider von Deimling, T., Ping, C.-L., Chanton, J. P., and Grosse, G.: Methane  
577 emissions proportional to permafrost carbon thawed in Arctic lakes since the 1950s, *Nat. Geosci.*, 9, 679–682,  
578 <https://doi.org/10.1038/ngeo2795>, 2016.

579 Watson, C. S., Quincey, D. J., Carrivick, J. L., and Smith, M. W.: The dynamics of supraglacial ponds in the Everest region,  
580 central Himalaya, *Glob. Planet. Change*, 142, 14–27, 2016.

581 Winslow, L. A., Read, J. S., Hanson, P. C., and Stanley, E. H.: Lake shoreline in the contiguous United States: quantity,  
582 distribution and sensitivity to observation resolution, *Freshw. Biol.*, 59, 213–223, 2014.

583 Xiong, G., Wang, G., Wang, D., Yang, W., Chen, Y., & Chen, Z. (2017). Spatio-temporal distribution of total nitrogen and  
584 phosphorus in Dianshan lake, China: The external loading and self-purification capability. *Sustainability*, 9(4), 500.

585 Yvon-Durocher, G., Hulatt, C. J., Woodward, G., and Trimmer, M.: Long-term warming amplifies shifts in the carbon cycle  
586 of experimental ponds, *Nat. Clim. Change*, 7, 209–213, <https://doi.org/10.1038/nclimate3229>, 2017.

587