



# BCUB - A large sample ungauged basin attribute dataset for British Columbia, Canada.

Daniel Kovacek<sup>1</sup> and Steven Weijs<sup>1</sup>

<sup>1</sup>Department of Civil Engineering, University of British Columbia, 2002 - 6250 Applied Science Lane, Vancouver, V6T 1Z4, BC, Canada

**Correspondence:** Daniel Kovacek (dkovacek@mail.ubc.ca)

**Abstract.** The British Columbia Ungauged Basin (BCUB) dataset is an open-source, extensible dataset of attributes describing terrain, soil, land cover, and climate indices of over one million ungauged basins in British Columbia, Canada including trans-boundary regions. The basin attributes included in the dataset follow those found in the large sample hydrology literature for their association with hydrological processes. The BCUB database is intended to support water resources research and practice, namely monitoring network analysis studies, or hydrological modelling where basin characterization is used for model calibration. The dataset, and the complete workflow to collect and process input data, to derive stream networks, delineate basins, and to extract basin attributes is available under a Creative Commons BY 4.0 license. The DOI link for the BCUB dataset is <https://doi.org/10.5683/SP3/JNKZVT> (Kovacek and Weijs, 2023).

## 1 Introduction

10 Spatial datasets available for geoscience research and practice are increasing in size, scale, resolution, and variety. Advances in the capture and processing of remote sensing data have in recent years led to open-access publication of continental and global scale geospatial datasets at high resolution (U.S. Geological Survey, 2022; Gleeson, 2018; Latifovic et al., 2010; Lehner et al., 2021; Thornton et al., 2021). We are well into the age of high quality, open-access geospatial data anticipated by Hrachowitz et al. (2013) following the decade of prediction in ungauged basins (PUB).

15 By contrast, the streamflow monitoring network in Canada has contracted over the last three decades. Based on the HYDAT dataset accessed at Environment Canada's national water data archive, the number of streamflow observation locations across Canada peaked in the order of 2300 in the 1980s, and reduced to roughly 1700 in 2022 (on average per day). According to surface water monitoring density standards developed by the World Meteorological Organization (WMO) (via Coulibaly et al. (2013)), nearly 90% of Canada's terrestrial area is under-monitored, and almost 40% is classified as ungauged. In general this trend holds for the province of British Columbia (BC), where outside of a few small regions in the south it is predominantly classified as ungauged or poorly gauged (Coulibaly et al., 2013).

The streamflow data used in a wide range of research and practice today comes from monitoring networks built over many decades, highlighting the significant lag between research aims of the present and monitoring layout decisions of the past. Monitoring network decisions today must anticipate information needs decades into the future.



25 Recent deep learning (DL) approaches to regional modeling use large sample datasets to infer relationships between climate  
input forcings and streamflow, and model performance improves when training incorporates static basin attributes (Kratzert  
et al., 2019). DL models benefit from training datasets (streamflow monitoring networks) representing basins that are diverse  
in geographic, hydrologic, and geophysical attributes, yet there is no clear consensus on how to evaluate networks in terms  
of diversity of attributes (Gauch et al., 2021). Increasing monitoring network diversity may be as simple as expanding the  
30 monitoring network according to the uniqueness of place described by Beven (2000), or a different approach can be to define  
a much larger set of ungauged basins and their hydrologically relevant attributes to use as a basis of comparison.

The vast and growing amount of geospatial information generated today requires considerable data assimilation effort to  
support specific research questions. A large, catchment-based dataset of geophysical attributes could support other disciplines  
that use basin attributes at the catchment level, for example in understanding changing water temperature and its effect on  
35 fish habitat (Daigle et al., 2017), or likewise for water quality monitoring in evaluating human-induced concentrations of toxic  
contaminants in fish (Scholes et al., 2016).

Water resource management decisions are typically made at the catchment level, so research and practice may be well served  
by datasets that are catchment-based, diverse in characteristics, and large in size and scale to reflect the scale-dependency of  
physical processes governing the rainfall-runoff response (Arsenault et al., 2020).

## 40 1.1 Motivation

The monitoring deficit of a region can be addressed by simply adding more stations, or under resource constraints optimal  
network arrangements can be approximated based on models trained on existing streamflow monitoring records, combined  
with information about unmonitored locations (Mishra and Coulibaly, 2010; Werstuck and Coulibaly, 2017, 2018). If large  
sample datasets improve predictability in ungauged locations by learning from diversity (Addor et al., 2017), a basis is needed  
45 to compare the existing monitoring network against the greater region it is intended to represent in relevant hydrological terms.  
The British Columbia Ungauged Basins (BCUB) (Kovacek and Weijs, 2023) is designed to be a dataset which i) uses only  
open access data sources that are continuous and complete over the study region, ii) is derived from the highest resolution DEM  
available to include smaller basins left out of other datasets, iii) is published under an open-source license, iv) is extensible  
both spatially and dimensionally to enable integration of new information as it is published, and v) is published with the full  
50 replication code based on widely used open-source libraries. Several existing datasets were reviewed for the desired qualities  
listed above, and for their potential to support research in network optimization, prediction in ungauged basins, and water  
resources more generally.

The BC Freshwater Atlas (FWA) (Gray, 2010) is the definitive source of freshwater feature mapping for British Columbia  
(BC). It contains roughly 3 million polygons representing the province-wide set of 1<sup>st</sup> order fundamental component water-  
shed units with a reference system designed to facilitate aggregation into larger watershed assessment units. The FWA dataset  
55 is strictly limited to the administrative bounds of BC, cutting off many important trans-boundary basins at borders. Since the  
dataset is primarily hydrographic, it does not include static basin attribute information commonly used in rainfall-runoff mod-



elling. The FWA is provided with an open-use license, but the code used to derive the dataset is to our knowledge unpublished, and as such it isn't readily replicable or extensible with consistent input data and methodology.

60 The National Hydrographic Network (NHN) (Geobase, 2004) contains a hydrographic feature set similar to the BC FWA. It covers all of Canada and includes trans-boundary basins along the US border, but the geometries are organized in Work Unit Limits (WULs) which do not represent complete basins. The watershed attributes are similarly limited, and the code used to derive the geometries is to our knowledge unpublished.

HydroSHEDS is a dataset for global-scale applications featuring river networks, watershed boundaries and other hydro-  
65 logical features derived from the NASA Shuttle Radar Topography Mission (SRTM) DEM for most of North America at a resolution of roughly 90m. At latitudes  $> 60^\circ$  North, corresponding to the northern border of BC with the Yukon territory, HydroSHEDS basins are derived from more coarse ( $\approx 500m$ ) Hydro1k (Wickel et al., 2007) elevation data. Attributes derived from distinct elevation data sources are difficult to compare as discussed in subsection 2.2, as the stream networks (and catchment boundaries) are unique to a DEM source and to the data processing methodology (Datta et al., 2022). Studies using the  
70 HydroSHEDS dataset typically exclude basins smaller than  $100 \text{ km}^2$  (Guth, 2011; Zhang et al., 2020; Kratzert et al., 2023).

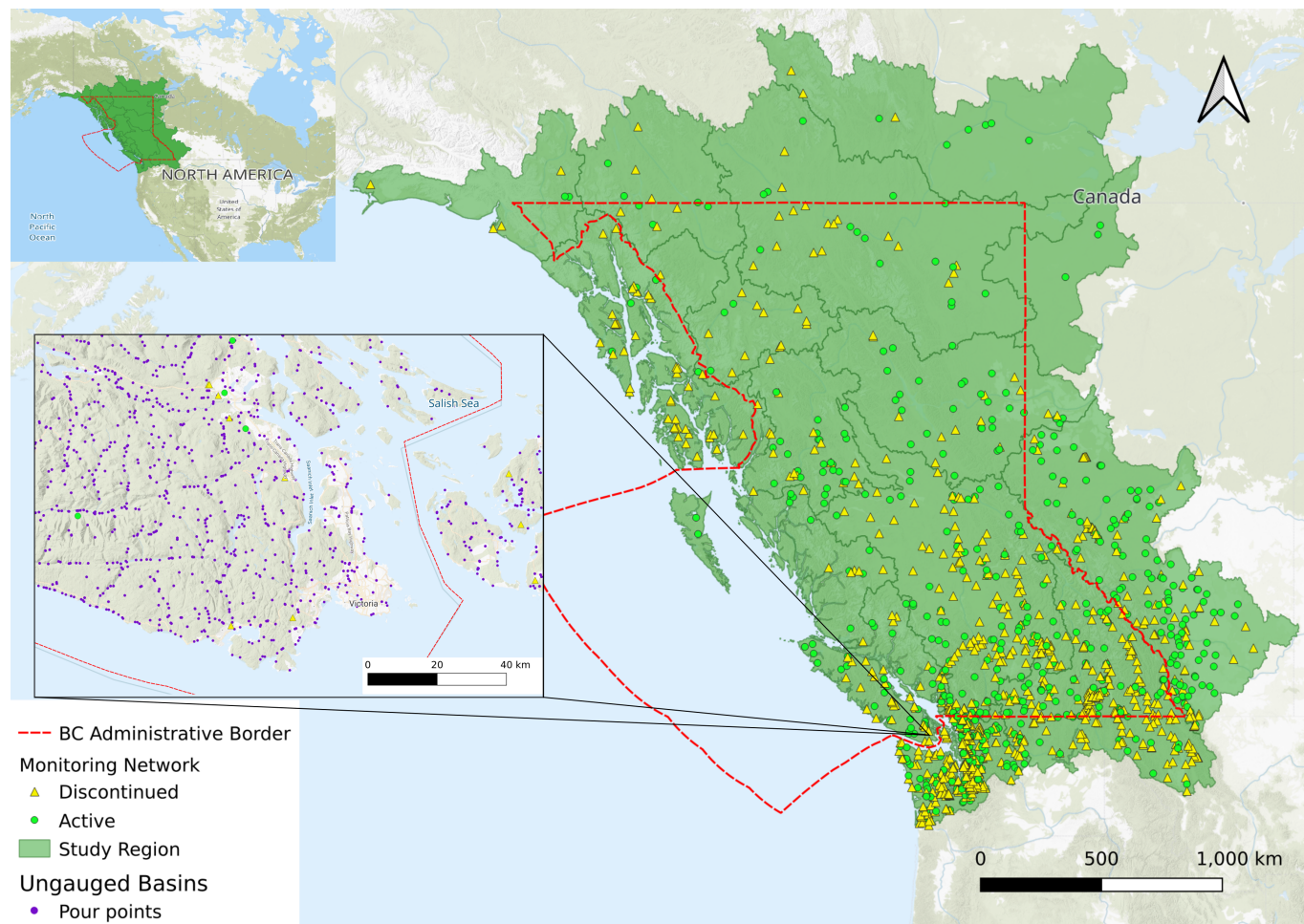
Large sample hydrology (LSH) datasets typically specify lower bounds on catchment area to filter out small basins due to uncertainty in basin delineation (Arsenault et al., 2020), and to ensure parameters are derived from sufficiently large samples (Guth, 2011), though explicit justification of a particular threshold is generally not provided. The HYSETS dataset (Arsenault et al., 2020) includes a caveat for attributes describing basins smaller than  $50 \text{ km}^2$ , representing nearly one third of the dataset.  
75 The uncertainty associated with such a large segment of the dataset (and the monitoring network it represents) highlights a gap that can be addressed in part with continuous and complete DEM coverage at greater resolution. The accuracy of stream network delineation improves with increasing DEM resolution (Tarolli and Dalla Fontana, 2009).

A large and diverse set of ungauged locations and associated attributes is sought to represent the decision space for network analysis and optimization, and more generally to support water resources research where catchment-based geospatial attributes  
80 are relevant.

## 1.2 British Columbia Ungauged Basin (BCUB) Database

The BCUB database contains a wide array of attributes describing terrain, land cover, soil permeability and porosity, and climate of over 1.2 million basins. Figure 1 shows the ungauged basin pour points representing the BCUB dataset, and the streamflow monitoring stations from the HYSETS dataset (Arsenault et al., 2020) that fall within the study region. The study  
85 region represents any terrestrial area within or upstream of any point within the BC administrative boundary (red dashed line in Figure 1), plus a buffer to include trans-boundary basins and to mitigate the edge selection bias of optimal sensor placement in random fields (Hershfield, 1965; Rouhani, 1985; Krause et al., 2006).

The attribute set describing each basin follows the HYSETS dataset as much as possible and includes select additional climate indices following the Camels dataset (Addor et al., 2017) to demonstrate how derived parameters can be added to the  
90 dataset. Three sets of land cover indices from the North American Land Change Monitoring System (NALCMS) (Latifovic



**Figure 1.** The study region (right) expands beyond the British Columbia administrative border to capture trans-boundary basin regions. Active and discontinued streamflow monitoring stations (those included in Arsenault et al. (2020)) are sparse and unevenly distributed as shown in the main figure at right, and the detail inset shows a sample of the high density of pour points defining catchments in the BCUB. (basemap from © MapTiler © OpenStreetMap contributors)

et al., 2010) associated with 2010, 2015, and 2020 are included to facilitate evaluation of land cover change at the basin level as called for by Addor et al. (2020). An example plot showing forest cover change between 2010 and 2020 is shown in section 3.

Following Wilkinson et al. (2016), to support knowledge discovery, innovation, and integration of data and methods in subsequent work, both the data and the code used to generate the data are openly available. The code is provided not to champion a particular method, but to highlight the nuance involved in developing large sample datasets that for brevity and clarity are typically left out of dataset description papers. There are no stochastic elements in the methodology, yet there are a large number of methodological choices that yield distinct outcomes. Providing the complete code at minimum aims to be explicit about these choices.



**Table 1.** Summary of catchment attribute source data.

Dataset	Attributes	Source
USGS 3DEP <sup>1</sup>	<b>Terrain:</b> area, elevation, aspect, slope	(U.S. Geological Survey, 2022)
GLHYMPS <sup>3</sup>	<b>Soil:</b> porosity, permeability	(Gleeson, 2018)
NALCMS <sup>2</sup>	<b>Land cover (2010, 2015, 2020):</b> forest, shrubs, grassland, wetland, crops, urban, water, snow and ice	(Latifovic et al., 2010)
DAYMET <sup>4</sup>	<b>Climate (daily estimates, 1980-2022):</b> precipitation, temperature, snow water equivalent, vapour pressure, shortwave radiation	(Thornton et al., 2022)

1. 3DEP: 3D Elevation Program, U.S. Geological Survey,

2. NALCMS: North American Land Change Monitoring System, accessed at <http://www.ccc.org/north-american-land-change-monitoring-system/>

3. Global Hydrogeology Maps.

4. Gridded daily climate estimates on a 1-km Grid for North America, Version 4. <https://daymet.ornl.gov/>

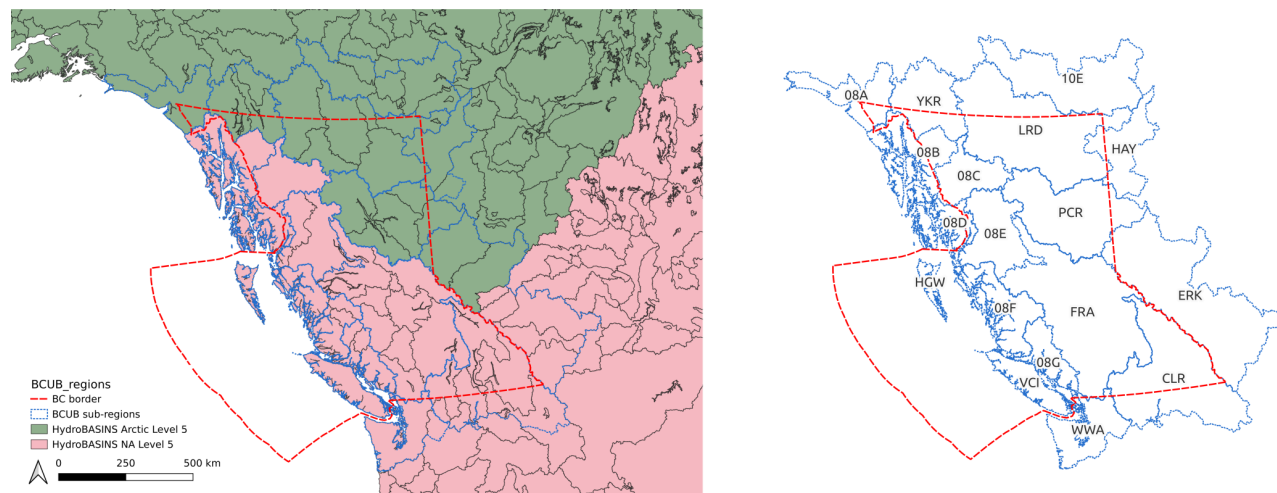
## 2 Data & Methods

### 100 2.1 Data collection and pre-processing

Static basin attributes in the BCUB dataset were extracted from the digital elevation, land cover, and soil geospatial layers described in Table 1 using basin polygons as clipping masks. Basin polygons were delineated from the set of pour points in the stream network representing river confluences. The stream network was derived from the 1 arc-second (30m at the equator) resolution USGS 3DEP (U.S. Geological Survey, 2022) digital elevation model (DEM) using the open-source software library  
105 Whitebox (version 2.3) (Lindsay, 2016).

The study region was divided into complete basin sub-regions as shown in Figure 2 (right) assembled from HydroBASINS (Lehner et al., 2021) data to simplify the automated basin delineation and attribute extraction work flow. Additional details about the data collection and pre-processing steps for generating the BCUB basin polygons are provided below.

- 110 1. **Hydrographic data:** Level 5 and 6 watersheds from the HydroBASINS dataset were used to break the study region into smaller complete basin sub-regions for data pre-processing. Lake polygons from HydroBASINS were used to filter out pour points in lakes that are vestiges of the stream network generation process.
2. **Digital elevation data:** The study region envelope, derived from the extents of the above geometries, was used to download the covering set of digital elevation tiles from the USGS 3D Elevation Program (U.S. Geological Survey, 2022). In addition, lower resolution (90m) DEM tiles from EarthEnv DEM90 (Robinson et al., 2014) were used in the  
115 data validation analysis presented in subsection 2.2.
3. **DEM raster processing:** Hydraulic conditioning of the DEM, including depression filling, resolving flats, computing flow direction and accumulation, and basin delineation were all done using the open-source geospatial analysis software Whitebox (Lindsay, 2016).



**Figure 2.** At right, the study region is divided into complete watershed sub-regions (encoded in the "region\_code" parameter) by merging level 5 & 6 HydroBASINS polygons to cover the BC boundary and include trans-boundary basins and include a minimum buffer of  $\approx 100$  km. The purpose of merging complete watershed regions is to manage computational resources the DEM pre-processing  $\rightarrow$  basin delineation  $\rightarrow$  attribute extraction pipeline.

### 2.1.1 Pour point set selection

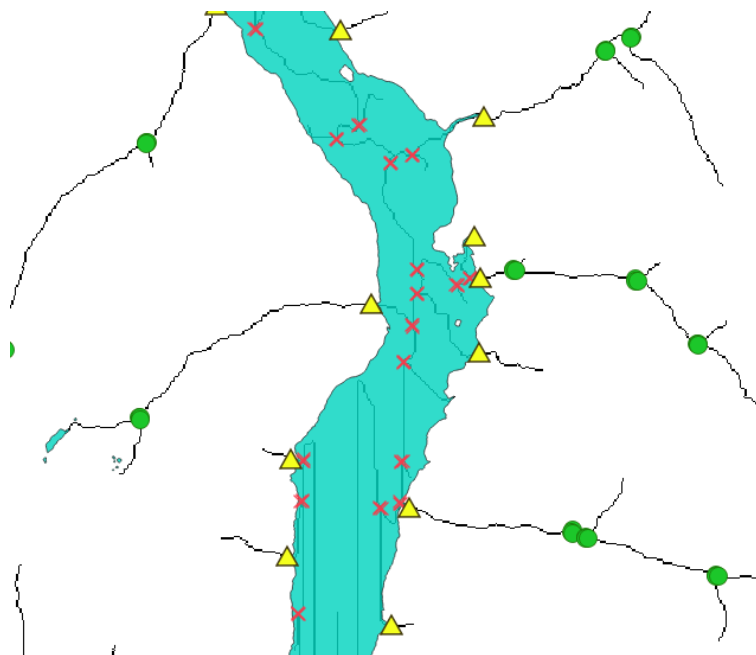
120 The basins in the BCUB database are delineated from a subset of raster cells representing the stream network. The set of  
pour points points used for basin delineation is called the *candidate monitoring location* (CML) set. By limiting the CML  
set to river confluences, the number of basins to process reduces to 5% of the complete set of stream network cells. Since  
changes in upstream accumulated area, and by extension hydrologic properties of the basin, are small along reaches between  
confluences, eliminating these points reduces redundancy and data processing, and reduces the decision space for subsequent  
125 network optimization analysis.

The CML set is defined by the following selection criteria:

1. **Confluences:** stream cells with more than two neighbouring stream cells, where the flow direction of more than one neighbouring stream cell is pointed toward the target cell, and
2. **River outlets:** intersections of river network lines with with ocean coastline, major basin outlets at the study region  
130 boundary, and intersections with lakes where the upstream contributing area is at least  $1 \text{ km}^2$ .

Stream confluences within lakes were excluded from the pour point set, as illustrated in Figure 3 where a red "x" denotes a spurious confluence within a lake, a yellow triangle represents the location where a river drains into a lake. Green circles represent upstream branches and their combination individually.

The stream network is defined by raster cells with a minimum upstream accumulation of  $1 \text{ km}^2$ . The headwaters mapped  
135 in the stream network raster are simply a vestige of the minimum area threshold used to define a stream network, so they



**Figure 3.** Example of river confluence (green circles) where spurious confluence points within lakes are excluded (red "x") and river-lake confluences are added (yellow triangle).

are not included in the pour point set. Accurate headwater identification (network extent mapping) requires a more rigorous approach to address uncertainty related to stream permanence (Shavers and Stanislawski, 2020). Mutzner et al. (2016) found classical (i.e. cumulative drainage area) threshold approaches do not capture spatial variability of headwater drainage networks in mountainous regions compared to detailed field survey mapping, and statistical methods are likewise unable to resolve local  
140 topography to accurately map headwater streams at low-resolution. Further discussion of uncertainty in stream networks is provided in subsection 2.2.

### 2.1.2 Attribute extraction

A basin polygon was delineated for each pour point in the CML set using the "unnest basins" function in the Whitebox software library (Lindsay, 2016). Basin attributes were derived for each basin by i) using basin polygons as raster clipping masks, and ii)  
145 spatial intersection of the basin polygon and geospatial raster and vector data in PostGIS (PostGIS Project Steering Committee and others, 2018). Basin attribute descriptions are provided in Table 2, and metadata attributes are described separately in Table 3. Attribute values are computed using the geometric mean of the raster pixel values contained in basin polygons in the case of soil permeability, the circular mean in the case of slope aspect, the fraction of total area in the case of land use, and the spatial mean for all other attributes. Additional pre-processing steps for the Daymet climate data are described in Table 2.



150 Two binary attributes are included in the attribute set. A soil flag value of 1 indicates that the sum of soil polygon areas clipped with the basin mask differs from the polygon area by more than 10%, to reflect where gaps exist in the GLHYMPS vector set. A permafrost flag value of 1 represents the presence of permafrost in the basin.

Expansion of the study region or addition of new attributes can be accomplished by following the processing methodology in the code repository provided. Four parameters derived from the Daymet daily precipitation data are processed in the code  
155 provided do demonstrate how computed parameters can be added from existing input data. The examples follow the Camels dataset Addor et al. (2017) and include:

1. **Low precipitation frequency:** frequency of days where precipitation  $< 1 \text{ mm day}^{-1}$ ,
2. **Low precipitation duration:** average duration of low precipitation events, or the number of consecutive low precipitation days  $< 1 \text{ mm day}^{-1}$ ,
- 160 3. **High precipitation frequency:** frequency of days where precipitation is  $\geq 5$  times the mean daily precipitation, and
4. **High precipitation duration:** average duration of consecutive high duration events, number of consecutive high precipitation days  $\geq 5$  times mean daily precipitation.

### 2.1.3 Data processing

Beyond data sources, the offline approach of deriving basins from source data and writing code to process attributes was  
165 adopted despite the elegant online polygon aggregation and processing approach demonstrated by Kratzert et al. (2023) in developing the Caravan dataset with use of Google Earth Engine (GEE) (Gorelick et al., 2017). Such an approach is preferable from the perspective of standardized methods of basin attribute extraction, but for our target of ungauged basins it does not eliminate the need for DEM pre-processing to generate stream networks, for filtering and extracting pour points, or for basin delineation. These steps represent a substantial portion of the basin attribute extraction workflow, and the what remains to  
170 process with GEE is still subject to usage limits, namely for processing the very large set of polygons, even considering an aggregated polygon approach.

A benefit of the offline approach is generating set of basin polygons from higher resolution DEM that is continuous and complete, and ensuring basin polygons match the DEM source from which terrain attributes are derived.

## 2.2 Technical Validation

175 The large number of basins in the BCUB dataset requires an automated approach to stream network validation and the basin polygons used to capture basin attributes. The representativeness of basin attributes is a function of the accuracy of the stream network derived from DEM. Higher resolution DEM can better resolve lower-relief topographic features resulting in better basin delineation performance, particularly for small basins (Zhang and Montgomery, 1994; Tarolli and Dalla Fontana, 2009; Woodrow et al., 2016).

180 It is important to emphasize that the  $1 \text{ km}^2$  minimum basin area threshold introduces significant uncertainty in the accuracy of the smallest basins, and in basins where topographic relief is low. Detailed validation of stream network accuracy is left to





future work that the BCUB is intended to support, and validation of the smallest basins used in studies is left to the user. Next we discuss indirect attribute validation methods, and limitations of the dataset and methods. The code to replicate the figures in this section is provided in the associated Github repository in the "validation" folder.

### 185 2.2.1 Vestigial effects of DEM resolution

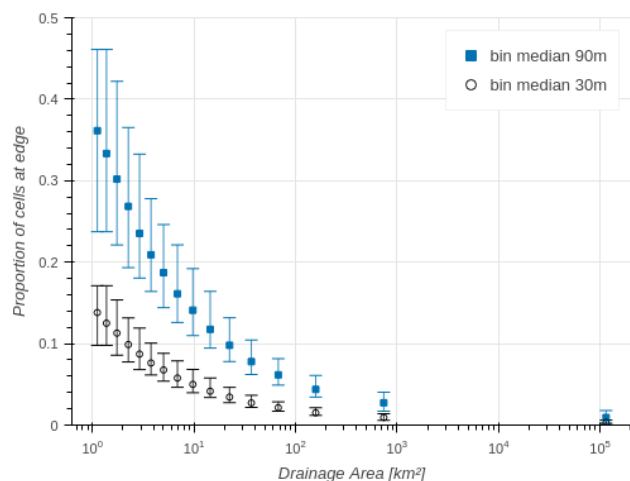
In addition to the process of hydraulic conditioning for stream network derivation, the grid representation of elevation introduces vestigial artifacts in the representation of basins and by extension in the extraction of basin attributes.

The stream network derived from DEM does not capture permanent water bodies, resulting in spurious river confluences. These vestigial confluences were excluded by using the lakes geometry layer from HydroBASINS as a mask, as described in  
190 subsection 2.1.1. Since HydroBASINS is derived from different sources, geometries do not align exactly with the stream network we derived from the 1 arc-second DEM.

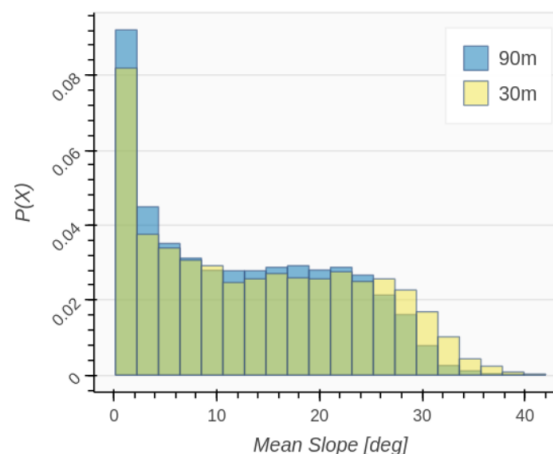
The disk space required to store a polygon is a linear function of the number of vertices defining it, and the precision of geographic coordinates describing the geometry. The basin polygons are simplified (using the Shapely library (Gillies, 2021) "simplify" function) using a tolerance equal to one-half the diagonal length of the raster pixel resolution. Simplifying (or  
195 smoothing) basin polygons is a trade-off between reducing the disk and bandwidth required to store and transmit large sets of basin geometries, and the representativeness of attributes that are captured by intersecting basin polygons with the various geospatial raster layers. The effect of polygon simplification is discussed in more detail in subsection 2.2.2.

The set of raster pixels representing each basin are captured using the "crop-to-cutline" function from the open source GDAL library (GDAL/OGR contributors, 2023) which by default captures pixels whose centroid lies within the polygon  
200 (pixels are not points, but quadrilaterals). Alternatively the larger set of *intersecting* pixels can be selected by setting the "CUTLINE\_ALL\_TOUCHED=TRUE" keyword argument. As basin area decreases (or raster resolution decreases), the difference in edge pixel selection method represents an increasing proportion of total pixels which may then yield significant differences in attribute values depending upon the clipping method used. Figure 4 shows the proportion of edge pixels representing the basin increases with decreasing area on a large sample of basin polygons and compare USGS 3DEP (30m at the equator) and  
205 EarthENV DEM90 (90m at the equator) DEM (EENV) on the same sample of polygons. The BCUB is derived from the USGS 3DEP (30m at the equator) DEM, and this exercise highlights one source of uncertainty introduced by the data processing methodology and suggests at what scale of basin the choice of clipping method becomes significant.

Mean basin slope is a widely used attribute (Addor et al., 2017; Alvarez-Garreton et al., 2018), defined in Arsenault et al. (2020) as "the average slope when considering the individual elevation differences between tiles" (raster pixels) and describes  
210 the basin's topographic relief. We used the WhiteboxTools Slope gradient function which computes slope for each DEM pixel using a 3rd-order Taylor polynomial fit (Florinsky, 2016) with a kernel size of 5x5 pixels. Mean basin slope increases with increasing resolution because topographic relief is better captured at higher resolution (Zhang and Montgomery, 1994). Figure 5 compares mean basin slope for two DEM sources with 30m and 90m resolution, where the higher resolution DEM is able to resolve greater topographic detail. The comparison is based on a random sample of 10K basins in the BCUB dataset  
215 ranging in size from 1 km<sup>2</sup> to 2 × 10<sup>5</sup> km<sup>2</sup>.



**Figure 4.** As the basin area decreases, the number edge pixels becomes a significant proportion of the total number of pixels representing the basin. Points in the above figure represent bin median values based on equiprobable binning ( $N \approx 600$  samples per bin), and the whiskers represent the 5 and 95 percentile values for each bin.



**Figure 5.** Higher resolution DEM captures greater topographic relief as shown by comparing the distribution of mean slope between 30m (USGS 3DEP) and 90m (EarthEnv) DEM on a random sample of 10,000 basins.

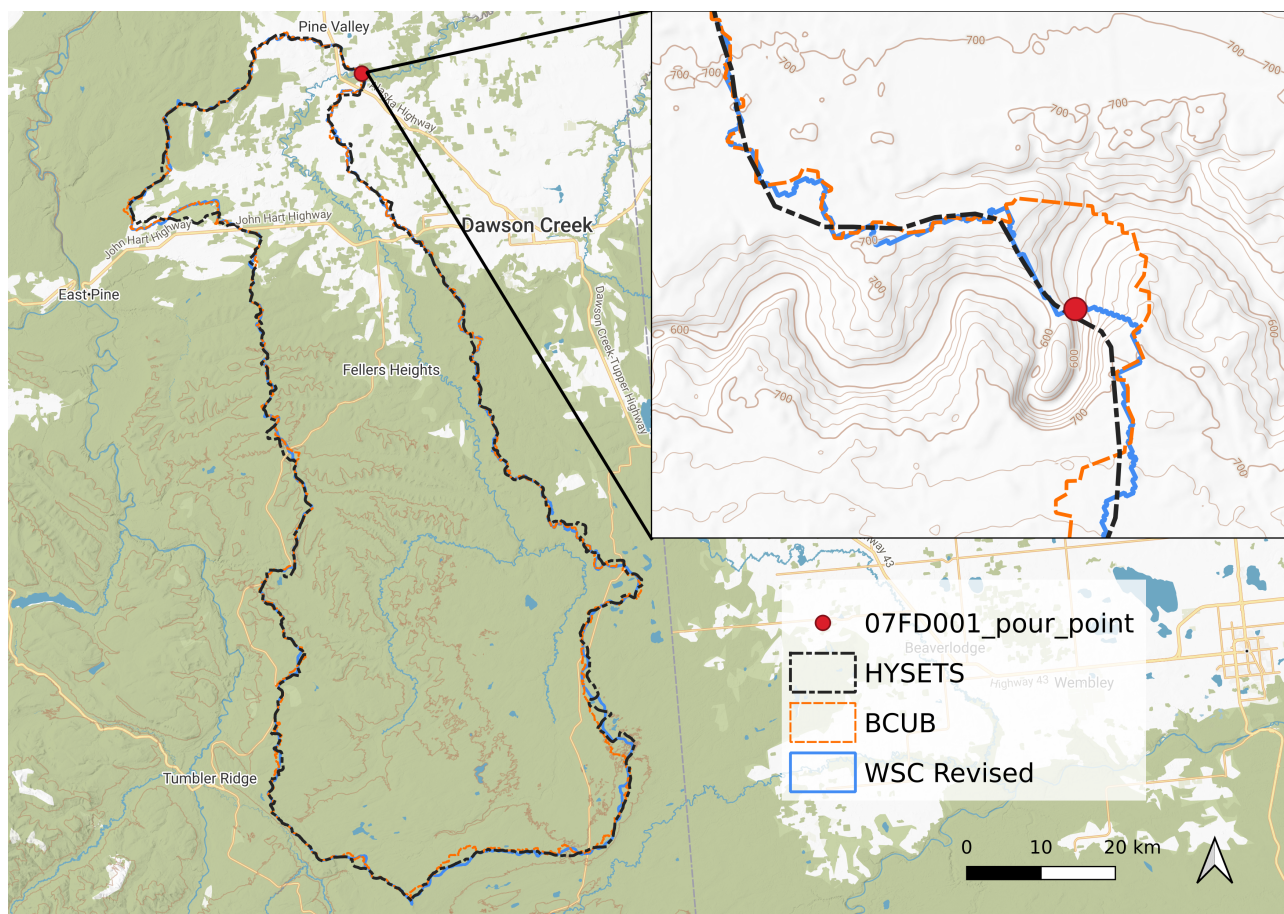
The sample of basins in Figure 5 shows a bias toward lower calculated mean slope from lower resolution DEM sources using the same basin polygon to capture pixels. Further interpretation of these differences is left to future work.

### 2.2.2 Basin Attributes and Self-Similarity

Mandelbrot (1967) described the measurement of coastline length as a function of the scale of observation, and the lines describing features like catchment boundaries and stream networks also exhibit self-similarity. Basin perimeter, stream gradient, and shape factors like elongation or compactness are length-based attributes used in many LSH datasets (Arsenault et al., 2020; Klingler et al., 2021; Kratzert et al., 2023). The basin compactness coefficient is defined as the ratio of basin perimeter to the circumference of a circle with equal basin area ((Gravelius, 1914) as cited in (Sassolas-Serrayet et al., 2018)). Length-based attributes are not comparable without consistent input DEM resolution and data pre-processing.

The difference in catchment boundary lines shown in Figure 6 illustrates why perimeter measurement can vary considerably due to input DEM resolution or basin delineation methodology. Basin perimeter is not included in the BCUB attribute set because unless otherwise treated, basin polygons derived from higher resolution DEM will measure longer perimeter.

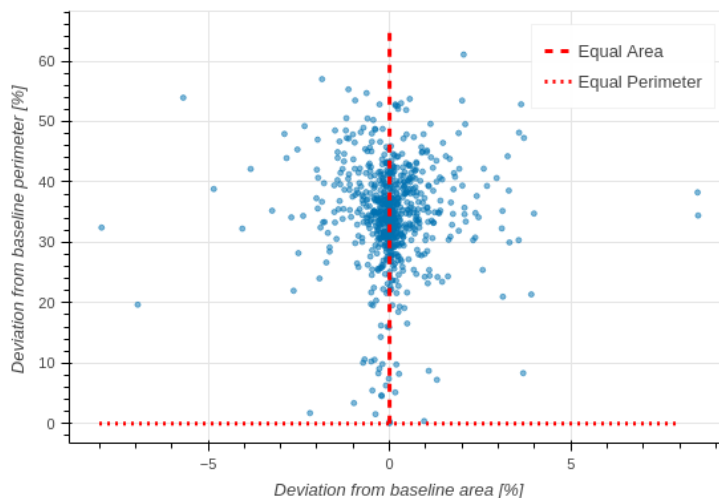
A large number of polygons used in LSH datasets (HYSETS and Caravan) were revised by the Water Survey of Canada in July 2022 and can be accessed at the (WSC) National Water Data Archive. We found all polygons common to both the HYSETS dataset and this updated polygon set, and computed pairwise comparisons of perimeter lengths. The sample used



**Figure 6.** An example edge detail of the same catchment boundary from three different sources where the intersecting area is over 98% of the published value. The HYSETS dataset polygon (back dash-dot line) comes from an earlier revision published by the WSC representing the Kiskatinaw River near Farmington (WSC ID 07FD001), while a recent revision (July 2022) by the WSC (solid blue) shows a distinct difference in polygon edges. The polygon from the BCUB (dashed orange) derived from USGS 3DEP DEM is different from both. (basemap from © MapTiler © OpenStreetMap contributors)

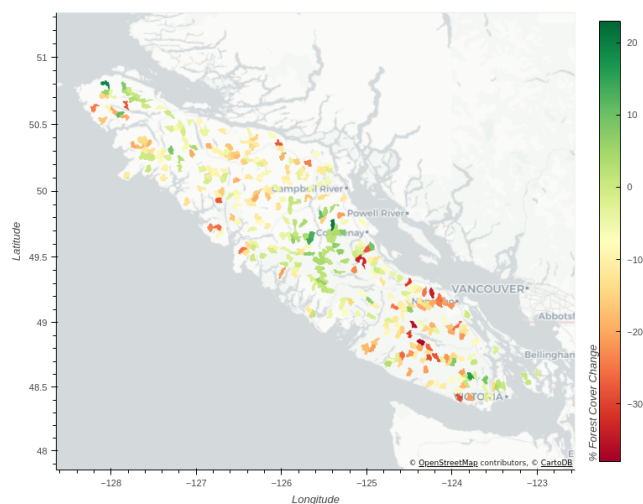
for the perimeter comparison includes 715 basins where the original and updated polygons were a close match. Similarity was evaluated based on the ratio of intersecting area to union area  $\geq 95\%$  to control for significant changes in the polygon shape. Figure 7 shows the newer revision polygon perimeter measurements are substantially greater, and that the deviation is not sensitive to the basin scale. This difference highlights the need to ensure consistent, continuous input DEM and data processing methodology if length-based attributes are included basin attribute datasets.

There were 1035 basin polygon revisions that did not meet the similarity criteria, reflecting the difficulty in retrospectively determining streamflow monitoring station locations from historical records (Arsenault et al., 2020).

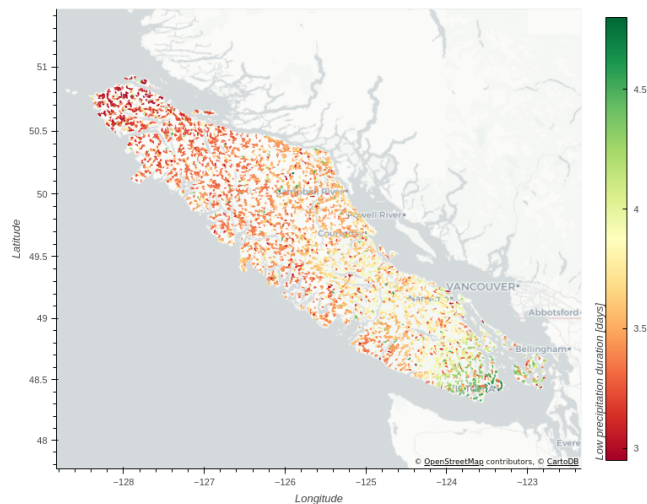


**Figure 7.** The DEM resolution and the processing methods used to derive basin polygons affects the measurement of perimeter. Polygons derived from different sources or using different methodologies will yield different values. Comparing sequential revisions of the same streamflow monitoring stations, the perimeter length is significantly different despite the area being roughly constant, and despite a close match between polygons according to a Jaccard Similarity Index match of  $\geq 95\%$ .

Average stream gradient is a length-based basin attribute that is a function of both raster resolution and the assumed location of channel head, usually by minimum area threshold. Robinson et al. (2014) calculated mean stream gradient as the ratio of the maximum total elevation change in the basin stream network to the length of the corresponding river reach. Stream length is a function of DEM resolution, and the length of reach is measured from the basin outlet to an uncertain headwater location (Hafen et al., 2020, 2022). In the derivation of the stream network for the BCUB dataset, headwater locations are simply a vestige of the assumed minimum drainage area threshold, and as a result an attribute representing average stream gradient is not included in the BCUB database.



**Figure 8.** An example visualization using the BCUB dataset maps the percent change in forest cover (as a percentage of the basin area) for basins with drainage area between 20 and 25 km<sup>2</sup> on Vancouver Island (VCI). Basemap from © OpenStreetMap contributors.



**Figure 9.** An example visualization using the BCUB dataset maps the mean annual duration of low precipitation (< 0.1 mm/day) for basins with drainage area between 2 and 5 km<sup>2</sup> on Vancouver Island (VCI). Basemap from © OpenStreetMap contributors.

### 245 3 Usage Notes

It is the hope that the BCUB dataset will serve a wide range of water resource research and practice where catchment-based attributes are integral to the methodology, or perhaps more importantly to express the limits of appropriate use and interpretation. Figure 8 and Figure 9 provide two basic examples of the kind of basin-level querying the BCUB is designed to support. Figure 8 shows basin-level changes in forest cover between 2010 and 2020 for basins in the range of 20 to 25 km<sup>2</sup>, and Figure 9 shows the mean duration of dry periods (days with less than 0.1 mm rainfall) for basins between 2 and 5 km<sup>2</sup>.

Stream networks are unique to the input DEM, and they are affected by the choice of pre-processing steps. The greatest degree of uncertainty is associated with the smallest basins with the lowest topographic relief. Zhang and Montgomery (1994) provides guidance about interpreting features at scales relative to DEM resolution. The representativeness of stream networks, and by extension basin polygons and the attributes captured by polygon masking, is an important component of uncertainty analysis and data reliability assessment. This aspect of the analysis is left to future work that the BCUB dataset is designed to support, in particular the lower limit of basin scale that can be supported by 1 arc-second DEM.

### 4 Code and data availability

The BCUB dataset (Kovacek and Weijs, 2023) is accessible under a Creative Commons BY 4.0 license through the Borealis data repository at <https://doi.org/10.5683/SP3/JNKZVT>. A summary of the dataset contents and supporting information is



260 presented in Table 4. The basin polygon geometries are provided in the open-source, cross-language Apache Parquet format (<https://parquet.apache.org/>), which has the convenience of supporting multiple geometries. The Parquet file format is supported by several widely used Python libraries, including Dask (<https://docs.dask.org/>) and GeoPandas (<https://geopandas.org/>), and the Arrow package features an interface for the R programming language (<https://arrow.apache.org/docs/r/>). The dask-geopandas library in Python (<https://dask-geopandas.readthedocs.io/>) is recommended for performance with large datasets.

265 The basin attributes are provided in two forms in the Borealis data repository. The larger form includes polygon geometries, and these are saved in the Parquet file format under the ‘basin\_polygons‘ folder (select the "tree" view to better organize the files for navigation). The Parquet file naming convention follows the sub-region codes shown in Figure 2. Sub-region geometries with their associated codes are provided for reference in BCUB\_regions\_4326.geojson. A "lighter" format is provided without basin geometries in comma delimited format in BCUB\_attributes\_20240117.csv. Referencing basins between files should be  
270 done by matching the pour point x,y coordinates included in both since these are by definition unique. Metadata describing the dataset is provided in MetaData.pdf, and additional basin attribute information, including descriptions and sources is provided in the Readme.pdf.

The scripts used to derive the dataset, and the validation results and figures shown in this paper are provided in an open-access Github repository (<https://github.com/dankovacek/bcub>). Figures 1 to 3 and 6 were prepared with the QGIS software  
275 (QGIS Development Team, 2023), and all remaining figures were created using the Bokeh data visualization library (Bokeh Development Team, 2023) in Python.

In addition, an example guide is provided through a set of Jupyter (Kluyver et al., 2016) notebooks ([https://dankovacek.github.io/bcub\\_demo/](https://dankovacek.github.io/bcub_demo/)) to demonstrate the complete process of data retrieval, pre-processing, basin delin-  
280 eation, attribute extraction, and data product usage. The code to produce Figure 8 using the Parquet file format is demonstrated in the final chapter of the Jupyter book demo, titled "Parquet Import and Usage".

*Author contributions.* Daniel Kovacek wrote the code to create the database and the Jupyter Notebook examples, and wrote the manuscript, and Steven Weijs provided research supervision and manuscript review.

*Competing interests.* The authors declare no competing interests.

*Acknowledgements.* This study received financial support from the British Columbia Ministry of Environment and Climate Change Strategy  
285 (Agreement #TP23EPEMA0031MY). The authors wish to express gratitude to all those contributing to open-source scientific software.



## References

- Addor, N., Newman, A. J., Mizukami, N., and Clark, M. P.: The CAMELS data set: catchment attributes and meteorology for large-sample studies, *Hydrology and Earth System Sciences*, 21, 5293–5313, 2017.
- Addor, N., Do, H. X., Alvarez-Garreton, C., Coxon, G., Fowler, K., and Mendoza, P. A.: Large-sample hydrology: recent progress, guidelines  
290 for new datasets and grand challenges, *Hydrological Sciences Journal*, 65, 712–725, 2020.
- Alvarez-Garreton, C., Mendoza, P. A., Boisier, J. P., Addor, N., Galleguillos, M., Zambrano-Bigiarini, M., Lara, A., Puelma, C., Cortes, G., Garreaud, R., et al.: The CAMELS-CL dataset: catchment attributes and meteorology for large sample studies–Chile dataset, *Hydrology and Earth System Sciences*, 22, 5817–5846, 2018.
- Arsenault, R., Brissette, F., Martel, J.-L., Troin, M., Lévesque, G., Davidson-Chaput, J., Gonzalez, M. C., Ameli, A., and Poulin, A.: A  
295 comprehensive, multisource database for hydrometeorological modeling of 14,425 North American watersheds, *Scientific Data*, 7, 1–12, 2020.
- Beven, K. J.: Uniqueness of place and process representations in hydrological modelling, *Hydrology and earth system sciences*, 4, 203–213, 2000.
- Bokeh Development Team: Bokeh: Python library for interactive visualization, <https://bokeh.org>, 2023.
- 300 Coulibaly, P., Samuel, J., Pietroniro, A., and Harvey, D.: Evaluation of Canadian National Hydrometric Network density based on WMO 2008 standards, *Canadian Water Resources Journal*, 38, 159–167, 2013.
- Daigle, A., Caudron, A., Vigier, L., and Pella, H.: Optimization methodology for a river temperature monitoring network for the characterization of fish thermal habitat, *Hydrological Sciences Journal*, 62, 483–497, 2017.
- Datta, S., Karmakar, S., Mezbahuddin, S., Hossain, M. M., Chaudhary, B. S., Hoque, M. E., Abdullah Al Mamun, M., and Baul, T. K.: The  
305 limits of watershed delineation: implications of different DEMs, DEM resolutions, and area threshold values, *Hydrology Research*, 53, 1047–1062, 2022.
- Florinsky, I.: *Digital terrain analysis in soil science and geology*, Academic Press, 2016.
- Gauch, M., Mai, J., and Lin, J.: The proper care and feeding of CAMELS: How limited training data affects streamflow prediction, *Environmental Modelling & Software*, 135, 104926, 2021.
- 310 GDAL/OGR contributors: GDAL/OGR Geospatial Data Abstraction software Library, Open Source Geospatial Foundation, <https://doi.org/10.5281/zenodo.5884351>, 2023.
- Geobase: National Hydro Network Data Production Catalogue, Available at: <https://www.nrcan.gc.ca/maps-tools-and-publications/survey-plans-and-data/standards-guidelines/10780>, accessed: 2023-04-30, 2004.
- Gillies, S.: Shapely: Manipulation and analysis of geometric objects, <https://github.com/Toblerity/Shapely>, 2021.
- 315 Gleeson, T.: GLObal HYdrogeology MaPS (GLHYMPS) of permeability and porosity, <https://doi.org/10.5683/SP2/DLGXYO>, 2018.
- Gorelick, N., Hancher, M., Dixon, M., Ilyushchenko, S., Thau, D., and Moore, R.: Google Earth Engine: Planetary-scale geospatial analysis for everyone, *Remote sensing of Environment*, 202, 18–27, 2017.
- Gravelius, H.: *Grundrifi der gesamten Gewisserkunde. Band I: Flufikunde (Compendium of Hydrology, vol. I. Rivers, in German)*, Goschen, Berlin, 1914.
- 320 Gray, M.: *Freshwater water atlas user guide*, GeoBC Integrated Land Management Bureau. Victoria, BC, 2010.
- Guth, P.: Drainage basin morphometry: a global snapshot from the shuttle radar topography mission, *Hydrology and Earth System Sciences*, 15, 2091–2099, 2011.



- Hafen, K. C., Blasch, K. W., Rea, A., Sando, R., and Gessler, P. E.: The influence of climate variability on the accuracy of NHD perennial and nonperennial stream classifications, *JAWRA Journal of the American Water Resources Association*, 56, 903–916, 2020.
- 325 Hafen, K. C., Blasch, K. W., Gessler, P. E., Sando, R., and Rea, A.: Precision of headwater stream permanence estimates from a monthly water balance model in the Pacific Northwest, USA, *Water*, 14, 895, 2022.
- Hershfield, D. M.: On the spacing of raingages, in: *Proceedings of the WMO/IASH Symposium on Design of Hydrometeorologic Networks*, Int. Assoc. Sci. Hydrol. Publ. vol. 67, pp. 72–79, Citeseer, 1965.
- Hrachowitz, M., Savenije, H., Blöschl, G., McDonnell, J., Sivapalan, M., Pomeroy, J., Arheimer, B., Blume, T., Clark, M., Ehret, U., et al.: A decade of Predictions in Ungauged Basins (PUB)—a review, *Hydrological sciences journal*, 58, 1198–1255, 2013.
- 330 Klingler, C., Schulz, K., and Herrnegger, M.: LamaH-CE: LARge-SaMple DATA for Hydrology and Environmental Sciences for Central Europe, *Earth System Science Data*, 13, 4529–4565, 2021.
- Kluyver, T., Ragan-Kelley, B., Pérez, F., Granger, B. E., Bussonnier, M., Frederic, J., Kelley, K., Hamrick, J. B., Grout, J., Corlay, S., et al.: Jupyter Notebooks—a publishing format for reproducible computational workflows., *Elpub*, 2016, 87–90, 2016.
- 335 Kovacek, D. and Weijis, S.: British Columbia Ungauged Basins Dataset, <https://doi.org/10.5683/SP3/JNKZVT>, 2023.
- Kratzert, F., Klotz, D., Herrnegger, M., Sampson, A. K., Hochreiter, S., and Nearing, G. S.: Toward improved predictions in ungauged basins: Exploiting the power of machine learning, *Water Resources Research*, 55, 11 344–11 354, 2019.
- Kratzert, F., Nearing, G., Addor, N., Erickson, T., Gauch, M., Gilon, O., Gudmundsson, L., Hassidim, A., Klotz, D., Nevo, S., et al.: Caravan—A global community dataset for large-sample hydrology, *Scientific Data*, 10, 61, 2023.
- 340 Krause, A., Guestrin, C., Gupta, A., and Kleinberg, J.: Near-optimal sensor placements: Maximizing information while minimizing communication cost, in: *Proceedings of the 5th international conference on Information processing in sensor networks*, pp. 2–10, 2006.
- Latifovic, R., Homer, C., Ressl, R., Pouliot, D., Hossain, S., Colditz Colditz, R., Olthof, I., Giri, C., and Victoria, A.: North American land change monitoring system (NALCMS), *Remote sensing of land use and land cover: principles and applications*. CRC Press, Boca Raton, 2010.
- 345 Lehner, B., Roth, A., Huber, M., Anand, M., Grill, G., Osterkamp, N., Tubbesing, R., Warmendinger, L., and Thieme, M.: HydroSHEDS v2. 0—Refined global river network and catchment delineations from TanDEM-X elevation data, in: *EGU General Assembly Conference Abstracts*, pp. EGU21–9277, 2021.
- Lindsay, J. B.: Whitebox GAT: A case study in geomorphometric analysis, *Computers & Geosciences*, 95, 75–84, 2016.
- Mandelbrot, B.: How long is the coast of Britain? Statistical self-similarity and fractional dimension, *science*, 156, 636–638, 1967.
- 350 Mishra, A. K. and Coulibaly, P.: Hydrometric network evaluation for Canadian watersheds, *Journal of Hydrology*, 380, 420–437, 2010.
- Mutzner, R., Tarolli, P., Sofia, G., Parlange, M. B., and Rinaldo, A.: Field study on drainage densities and rescaled width functions in a high-altitude alpine catchment, *Hydrological Processes*, 30, 2138–2152, 2016.
- PostGIS Project Steering Committee and others: PostGIS, spatial and geographic objects for PostgreSQL, <https://postgis.net>, 2018.
- QGIS Development Team: QGIS Geographic Information System, Open Source Geospatial Foundation, <http://qgis.org>, 2023.
- 355 Robinson, N., Regetz, J., and Guralnick, R. P.: EarthEnv-DEM90: A nearly-global, void-free, multi-scale smoothed, 90m digital elevation model from fused ASTER and SRTM data, *ISPRS Journal of Photogrammetry and Remote Sensing*, 87, 57–67, 2014.
- Rouhani, S.: Variance reduction analysis, *Water Resources Research*, 21, 837–846, 1985.
- Sassolas-Serrayet, T., Cattin, R., and Ferry, M.: The shape of watersheds, *Nature communications*, 9, 3791, 2018.
- Scholes, R. C., Hageman, K. J., Closs, G. P., Stirling, C. H., Reid, M. R., Gabriëlsson, R., and Augspurger, J. M.: Predictors of pesticide concentrations in freshwater trout—The role of life history, *Environmental Pollution*, 219, 253–261, 2016.
- 360





- Shavers, E. and Stanislawski, L. V.: Channel cross-section analysis for automated stream head identification, *Environmental Modelling & Software*, 132, 104 809, 2020.
- Tarolli, P. and Dalla Fontana, G.: Hillslope-to-valley transition morphology: New opportunities from high resolution DTMs, *Geomorphology*, 113, 47–56, 2009.
- 365 Thornton, M., Shrestha, R., Wei, Y., Thornton, P., Kao, S., and Wilson, B.: Daymet: Monthly Climate Summaries on a 1-km Grid for North America, Version 4 R1. ORNL DAAC, Oak Ridge, Tennessee, USA, 2022.
- Thornton, P. E., Shrestha, R., Thornton, M., Kao, S.-C., Wei, Y., and Wilson, B. E.: Gridded daily weather data for North America with comprehensive uncertainty quantification, *Scientific Data*, 8, 190, 2021.
- U.S. Geological Survey: 1 Arc-second Digital Elevation Models (DEMs) USGS National Map 3D Elevation Program, <https://data.usgs.gov/datacatalog/data/USGS:35f9c4d4-b113-4c8d-8691-47c428c29a5b>, [Online; accessed 3 March 2022], 2022.
- 370 Werstuck, C. and Coulibaly, P.: Hydrometric network design using dual entropy multi-objective optimization in the Ottawa River Basin, *Hydrology Research*, 48, 1639–1651, 2017.
- Werstuck, C. and Coulibaly, P.: Assessing Spatial Scale Effects on Hydrometric Network Design Using Entropy and Multi-objective Methods, *JAWRA Journal of the American Water Resources Association*, 54, 275–286, 2018.
- 375 Wickel, B., Lehner, B., and Sindorf, N.: HydroSHEDS: A global comprehensive hydrographic dataset, in: *AGU Fall Meeting Abstracts*, vol. 2007, pp. H11H–05, 2007.
- Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., Appleton, G., Axton, M., Baak, A., Blomberg, N., Boiten, J.-W., da Silva Santos, L. B., Bourne, P. E., et al.: The FAIR Guiding Principles for scientific data management and stewardship, *Scientific data*, 3, 1–9, 2016.
- Woodrow, K., Lindsay, J. B., and Berg, A. A.: Evaluating DEM conditioning techniques, elevation source data, and grid resolution for  
380 field-scale hydrological parameter extraction, *Journal of hydrology*, 540, 1022–1029, 2016.
- Zhang, J., Condon, L. E., Tran, H., and Maxwell, R. M.: A national topographic dataset for hydrological modeling over contiguous United States, *Earth System Science Data Discussions*, 2020, 1–26, 2020.
- Zhang, W. and Montgomery, D. R.: Digital elevation model grid size, landscape representation, and hydrologic simulations, *Water resources research*, 30, 1019–1028, 1994.



**Table 2.** Basin attributes in the BCUB database derived from USGS 3DEP (DEM), NALCMS (land cover), GLHYMPS (soil), and NASA Daymet (climate) datasets.

Group	Description (BCUB label)	Aggregation	Units
Terrain	Drainage Area (drainage_area_km2)	at pour point	$km^2$
	Elevation (elevation_m)	spatial mean	$m$ above sea level
	Terrain Slope (slope_deg)	spatial mean	$^\circ$ (degrees)
	Terrain Aspect (aspect_deg)	circular mean <sup>2</sup>	$^\circ$ (degrees)
Land Cover <sup>3</sup>	Cropland (land_use_crops_frac_<year>)		
	Forest (land_use_forest_frac_<year>)		
	Grassland (land_grass_forest_frac_<year>)		
	Shrubs (land_use_shrubs_frac_<year>)	spatial mean	% cover
	Snow & Ice (land_use_snow_ice_frac_<year>)		
	Urban (land_use_urban_frac_<year>)		
	Water (land_use_water_frac_<year>)		
Soil <sup>4</sup>	Wetland (land_use_wetland_frac_<year>)		
	Permeability (logk_ice_x100)	geometric mean	$m^2$
	Std. Dev. Permeability (k_stdev_x100)	geometric mean	$m^2$
Climate <sup>5</sup>	Porosity (porosity_x100)	spatial mean	% cover
	Precipitation (prcp)		mm
	Minimum Temperature (tmin)		Celsius
	Maximum Temperature (tmax)		Celsius
	Snow Water Equivalent (swe)		Celsius
	Shortwave Radiation (srad)	spatial and	$W/m^2$
	Vapour Pressure (vp)	temporal mean	Pa
	High precipitation frequency (high_prpc_freq)		days/year
	Low precipitation frequency (low_prpc_freq)		days/year
	High precipitation duration (high_prpc_duration)		days
Low precipitation duration (low_prpc_duration)		days	

1. Spatial aspect is expressed in degrees counter-clockwise from the east direction.

2. The <year> suffix specifies the land cover dataset (2010, 2015, or 2020),.

3. Soil parameters follow definitions from Gleeson (2018).

4. Only the climate parameters directly extracted from distinct daymet source variables are shown here. Additional computed parameters are discussed in subsection 2.1.3.

5. A high precipitation event is defined as total daily precipitation greater than 5x the annual mean, and the duration refers to the mean duration of high precipitation events.

6. A low precipitation event is defined as total daily precipitation less than 0.1mm, and the duration refers to the mean duration of low precipitation events.



**Table 3.** BCUB dataset metadata attributes.

Group	Description (BCUB label)	Aggregation	Units
	Region code identifier (region_code)	-	-
	Pour point <sup>1</sup> (ppt_x, ppt_y)	-	m
Metadata	Basin centroid <sup>1</sup> (centroid_x, centroid_y)	-	m
	Soil Flag (soil_flag)	-	binary (0/1)
	Permafrost Flag (permafrost_flag)	-	binary (0/1)

1. Geometries are projected to the BC Albers (EPSG:3005) coordinate reference system.

**Table 4.** Summary of data repository contents.

Filename	Description
<b>BCUB_attributes_20231114.csv</b>	Basin attributes with geographic coordinates describing the basin centroid and the basin outlet (pour point). Basin polygon geometries not included for performance.
<b>polygons/*.parquet</b>	Basin attributes and associated catchment boundary polygons are organized into sub-regions due to limit file sizes.
<b>BCUB_regions_4326.geojson</b>	Spatial reference file describing the study area sub-regions corresponding to parquet filename prefixes (i.e. VCI_basins.parquet)
<b>MetaData.pdf</b>	General information about the dataset content, formats, versioning, and input data sources.
<b>README.pdf</b>	Basin attribute descriptions and method references.