

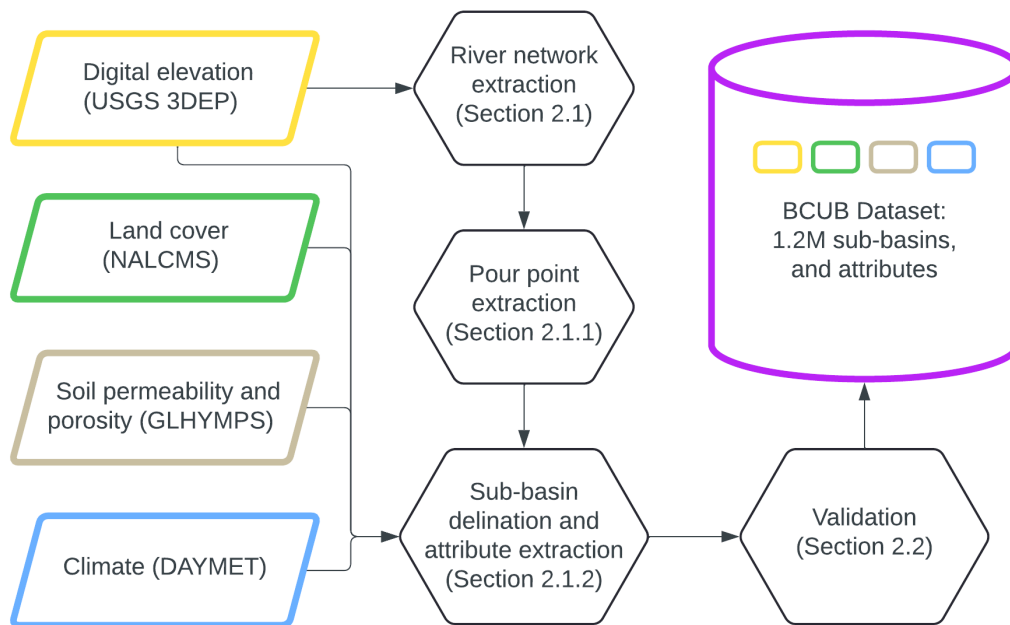
Author's Thanks (2024-05-21): We are grateful for the time taken and the effort made by the reviewer to consider our work and provide feedback. The points raised highlight important clarifications to both the content and delivery of the paper that undoubtedly improve its quality. Please see below for our responses to specific feedback.

RC1 Comment: A definition of the basin considered in this study is needed. Basin is a term that is interchangeable with catchment and watershed, but it typically refers to the entire drainage area of a river. In this article, 'basin' represents the local watershed of each river-reach. The term 'sub-catchment' or 'sub-basin' is more appropriate here.

Author's Response: As per your recommendation, we will add a definition of the term basin in our paper to clarify. While we agree that many of the basins considered in our dataset could be classified as sub-basins or sub-sub-basins, we use the term basin in a wider sense of the definition. This is in line with literature about ungauged basins. For example, the usage of "basin" in "A decade of Prediction in Ungauged Basins (PUB)--a review" (Hrachowitz et al., 2013) does not seem to refer exclusively to the entire drainage area of a river. To avoid confusion, we will explicitly define our use of "basin" at the start of Section 1.2.

RC1 Comment: To understand the process more easily, a flowchart showing different steps of BCUB database development in the methodology section would be helpful.

Author's Response: We agree a diagram will provide a useful overview of the full process. The diagram below will be added to the manuscript to represent the dataset development process:

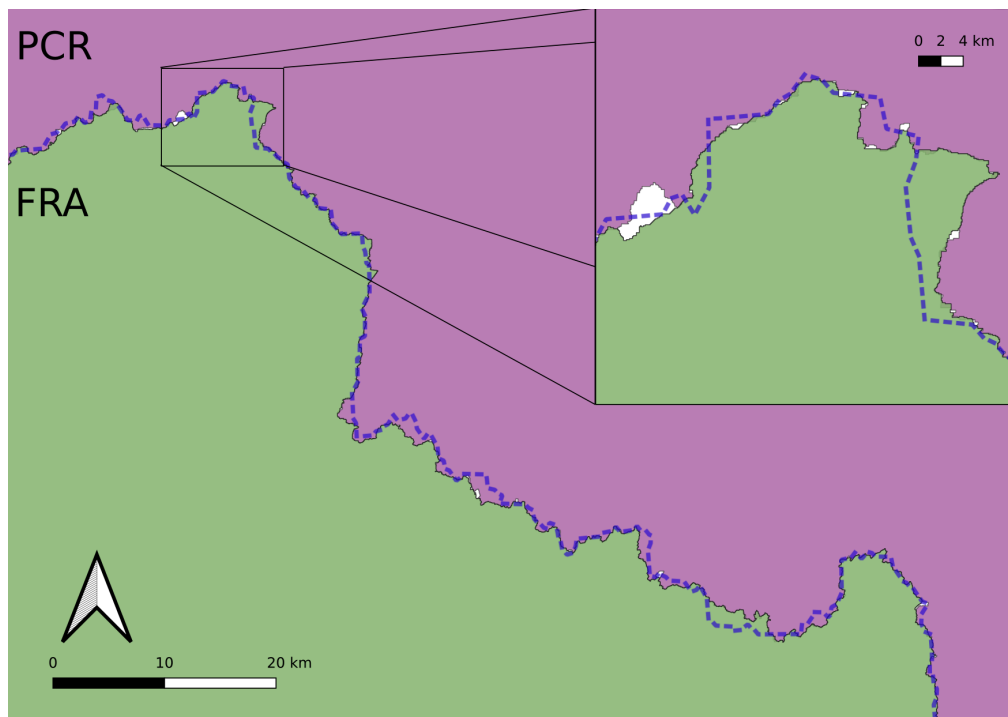


RC1 Comment: *The reason for using the HydroBASINS watersheds (level-5 and 6) to subdivide the study region is understandable. However, the underlying hydrography data in the HydroBASINS and BCUB databases are different. So, there is a chance of missing a part of the sub-catchments located near the regional boundary in the BCUB database. For example, a part of the sub-catchment of the PCR region, located near the boundary between PCR and FRA, may overshoot to the FRA region due to hydrographic data inconsistencies. How was this issue addressed during the development of this database?*

Author's Response: Thank you for raising this important point. While the region polygons assembled from HydroBASINS are a helpful tool for organizing the data processing pipeline, indeed their use yields different bounds whose effect on sub-basin delineation is in the order of the size of the smallest sub-basins in the BCUB dataset.

The Caravan dataset (Kratzert et al., 2023) clearly describes the issue with aggregating attributes from catchment polygons that do not align with the HydroBASINS dataset. By independently deriving the region bounds from a single continuous DEM source (USGS 3DEP 30m grid), we avoid the problem of misalignment with HydroBASINS polygons, however it does not solve the problem of disagreement in region bounds defined independently of HydroBASINS.

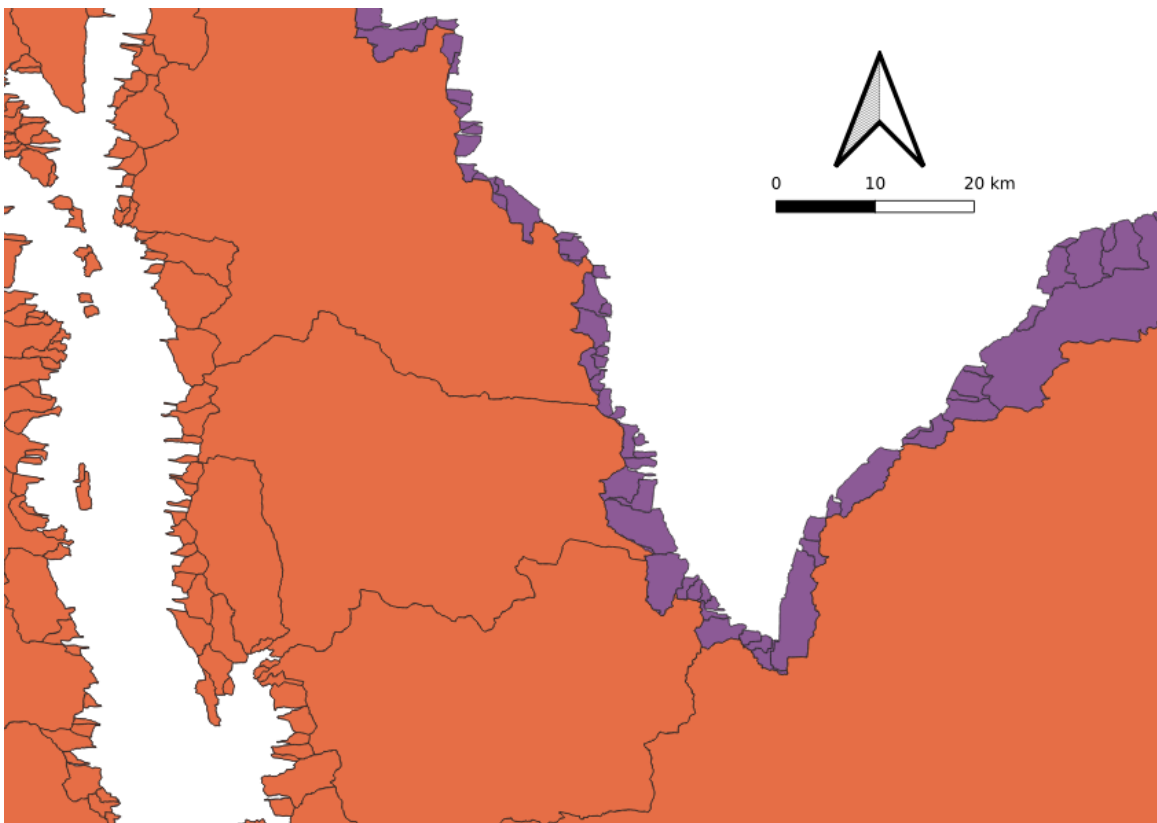
Below is an outline of the process we use to independently redefine sub-region polygons from the DEM and quantify uncertainty in region bounds in the BCUB dataset.



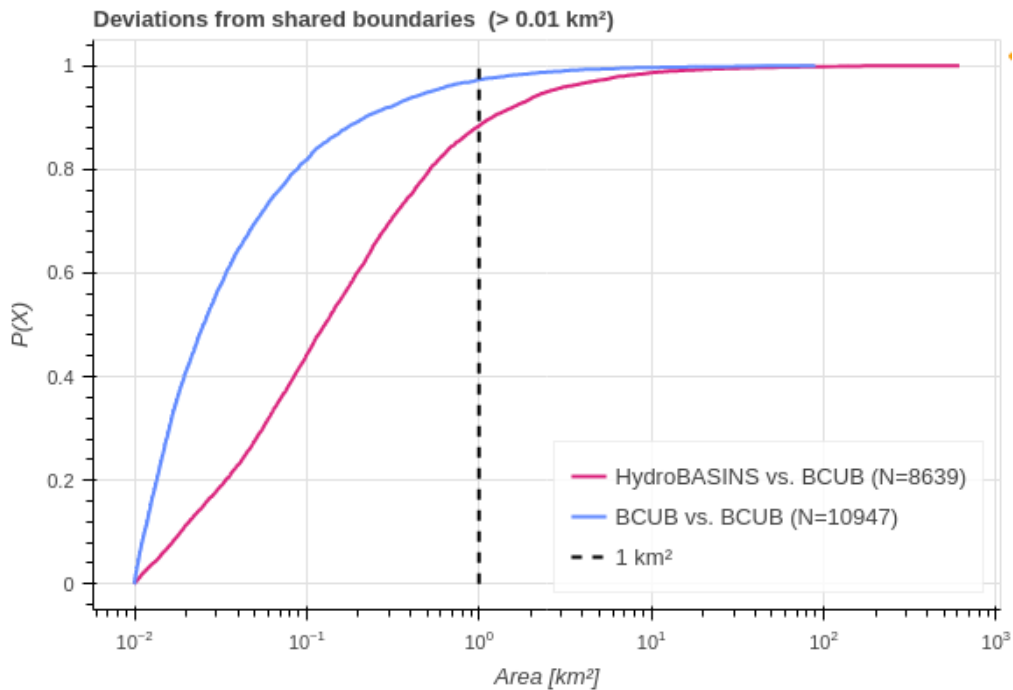
The edge detail inset in the figure above shows an example segment of region boundaries aggregated from HydroBASINS (blue dashed line) compared to independently derived region

bounds. The purple (Peace, PCR) and green (Fraser, FRA) coloured areas represent the region boundaries derived independently using the USGS 3DEP DEM (30m grid resolution), referred to here as the BCUB region boundaries. White areas are gaps that remain following the iterative boundary definition process described below. We define boundary deviations as polygons representing i) gaps between region bounds where the DEM resolution does not resolve which direction the small area drains, and ii) boundary overlaps when delineating from pour points in distinct basins with shared boundaries.

The process begins by applying a (5km) buffer to the region boundaries aggregated from level 5 and 6 HydroBASINS polygons, and using these buffered polygons as clipping masks on the DEM. The purpose of this step is to avoid restricting the catchment boundary delineation by the clipping mask. The covering set of polygons (catchments) are then delineated from the clipped DEM for each region, and the exterior edges (of the union of intersecting geometries) are checked to verify that they do not touch the edge of the buffered region polygon. Where the edges intersect, we manually expand the buffer (DEM clipping mask) in QGIS and re-derive the covering set of catchments until the buffer is sufficient, i.e. the covering set of basins does not touch the edge of the clipping mask. The use of a buffer causes small catchments to be delineated which drain to basins in adjacent regions, and these are excluded by identifying that they are completely contained by the clipping mask of the neighbouring region. The figure below illustrates the excluded vestigial edge sub-basins (purple) and the remaining covering basin set (orange).



Delineating region boundaries independently from the covering set of basins does not yield perfectly shared boundaries, but these deviations are substantially smaller than those resulting from aggregating the HydroBASINS levels 5 and 6 polygons. The distribution of the size of deviations from shared sub-region boundaries are shown in the figure below. The red series represents deviations between the BCUB region bounds and HydroBASINS-derived bounds (median area of 0.13 km²), while the blue series represents deviations (overlaps and gaps) within the BCUB sub-region boundaries (median area = 0.03 km²). Polygons smaller than 0.01 km², or 1% of the smallest sub-basin in the BCUB dataset were neglected.



We will incorporate a geometry flag attribute in the BCUB dataset for any sub-basin that intersects or touches at least one boundary deviation, and will include a decimal value to represent the total deviation area as a percentage of the sub-basin area. Where two different sub-basins claim the same area, either bordering sub-basin may overestimate the catchment area (indicated by a positive % value). Where an area is not claimed by any basin but is not necessarily endorheic, either bordering sub-basin may be underestimating the catchment area (indicated by a negative % value). The percentage represents the maximum expected percentage error from the uncertain boundary. The purpose of including these quantities is to communicate (some part of) the uncertainty in defining region bounds where the size of the uncertain area exceeds 1% of any sub-basin area. We will update the region boundaries in the data repository, and we will additionally provide the set of polygons representing boundary deviations as a .geojson file to facilitate corrections given updated information resolving these disagreements.

We additionally point out that a precise coastline definition (or ocean masking) at the resolution of the input DEM is important for the river network processing computation, otherwise vestigial river

segments occur in the ocean parallel to coastlines where the HydroBASINS polygons extend over ocean surface. We crop the coastline using the NALCMS land cover data ocean pixels – the land cover data are well suited to the input DEM since the both products are provided in the same grid resolution.

Finally, these region boundary updates will require revising the BCUB dataset. We will reprocess all affected sub-basins and update the dataset with the above additional information, namely the catchment delineation flag and the percent area represented by uncertain region boundaries. The additional detail provided here will appear in some form in the manuscript. The code used to derive the region boundary deviations will be provided along with the existing validation code in the open-source code repository. We believe these revisions will result in a more transparent and higher quality dataset, and we appreciate the reviewer raising this important detail.

***RC1 Comment:** It is sometimes difficult to follow the article due to inconsistencies in the statements. For example, the line 76 in the motivation section, “The accuracy of stream network delineation improves with increasing DEM resolution.” The transition from the previous lines to this one is not smooth.*

Author’s Response:

Agreed. This point is made in a more appropriate context later in the text (lines 180-185) so we have removed the statement.

***RC1 Comment:** Another example of inconsistency is in line 134, where the delineation of the stream network is discussed after the description of the pour point selection process from the stream network. It would be more appropriate to discuss the stream network delineation process before selecting pour points.*

Author’s Response:

Agreed. The order of stream network extraction and pour point selection have been adjusted accordingly to improve the consistency overall narrative and sequencing of arguments.

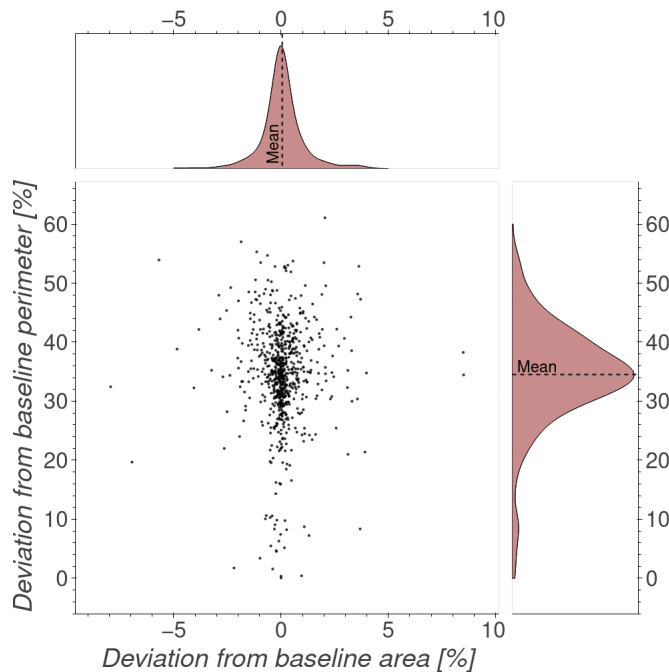
***RC1 Comment:** Line 103: Please provide the minimum drainage area threshold used to delineate the stream network from USGS 3DEP*

Author’s Response: The minimum drainage area threshold used is 1 km² which corresponds to the smallest sub-basin included in the HYSETS dataset (Arsenault et al. 2020) and to the smallest monitored basin in the British Columbia streamflow monitoring network. This reference is made explicit in the text, but your note identifies where (we agree) it should be placed earlier in the text. The text around line 103 has been updated to explicitly state the minimum threshold.

RC1 Comment: Figure 7: This is a nice figure to show the impact of using DEM with different resolutions. The plot with colored density would be more helpful to understand the figure.

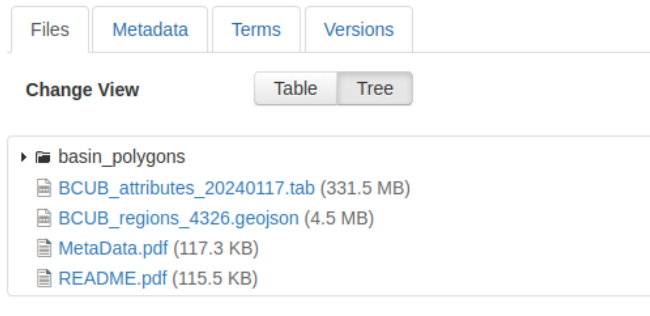
Author's Response:

Figure 7 has been modified to show the distributions in both x and y, as shown below, which we hope adds clarity to the meaning of the figure. We tried a 2D (kernel) density plot to unsatisfactory effect. We believe the addition of x and y distributions are a reasonable compromise to sufficiently demonstrate the point that when increasing the input DEM resolution, the mean change in area is near zero while the corresponding change in perimeter is substantially greater than zero. We also add that the coefficient of determination (R^2) between x and y (area and perimeter deviation from baseline) is 0.00, which means that the marginal distributions of x and y do not lose any information relative to the joint distribution of x,y. We will modify the figure as shown below and add the coefficient of determination to the manuscript highlighting this point.



RC1 Comment: When using QGIS version 3.28 to open the dataset, it displays the pour point location instead of the sub-basin polygon. Has the delineated sub-basin geometry been excluded from the database?

Author's Response: The tabular file (BCUB_attributes_20240117.tab) contains the x,y coordinates of the pour point (ppt) and basin centroid ('centroid_x', 'centroid_y', 'ppt_lon_m_3005', 'ppt_lat_m_3005') while due to the very large file sizes, the polygon geometries are provided separately in the Parquet file format saved under the "basin_polygons" folder in the data repository:



Parquet is supported by GDAL as of version 3.5, so QGIS must be compiled with GDAL \geq 3.5 which is not default in some environments.

Please see the following for information about versions and compatibility:

<https://gis.stackexchange.com/questions/430973/importing-geoparquet-file-in-qgis>

Reading/writing Parquet in R:

https://arrow.apache.org/docs/r/reference/read_parquet.html

Reading/writing Parquet in Python:

<https://arrow.apache.org/docs/python/parquet.html>

Parquet is also implemented in Julia, MATLAB, Rust, Go, Java, C++, and others:

<https://arrow.apache.org/docs/>

References:

1. Hrachowitz, Markus, et al. "A decade of Predictions in Ungauged Basins (PUB)—a review." *Hydrological sciences journal* 58.6 (2013): 1198-1255.
2. Arsenault, R., Brissette, F., Martel, J.-L., Troin, M., Lévesque, G., Davidson-Chaput, J., Gonzalez, M. C., Ameli, A., and Poulin, A.: "A comprehensive, multisource database for hydrometeorological modeling of 14,425 North American watersheds", *Scientific Data*, 7, 1-12, 2020.
3. Kratzert, Frederik, et al. "Caravan-A global community dataset for large-sample hydrology." *Scientific Data* 10.1 (2023): 61.