

# Deriving a Transformation Rate Map of Dissolved Organic Carbon over the Contiguous U.S.

Lingbo Li<sup>1</sup>, Hong-Yi Li<sup>1\*</sup>, Gita Abeshu<sup>2</sup>, Jinyun Tang<sup>3</sup>, L. Ruby Leung<sup>2</sup>, Chang Liao<sup>2</sup>, Zeli Tan<sup>2</sup>, Hanqin Tian<sup>4</sup>, Peter Thornton<sup>5</sup>, Xiaojuan Yang<sup>5</sup>

<sup>1</sup>Department of Civil and Environmental Engineering, University of Houston, Texas, USA

<sup>2</sup>Pacific Northwest National Laboratory, Washington, USA

<sup>3</sup>Lawrence Berkeley National Laboratory, California, USA

<sup>4</sup>Boston College, Massachusetts, USA

<sup>5</sup>Environmental Sciences Division, and Climate Change Science Institute, Oak Ridge National Laboratory, Tennessee, USA

Correspondence to: Hong-Yi Li ([hongyili.jadison@gmail.com](mailto:hongyili.jadison@gmail.com))

**Abstract.** Riverine dissolved organic carbon (DOC) plays a vital role in regional and global carbon cycles. However, the processes of DOC conversion from soil organic carbon (SOC) and leaching into rivers are insufficiently understood, inconsistently represented, and poorly parameterized, particularly in land surface and [earthEarth](#) system models. As a first attempt to fill this gap, we propose a generic formula that directly connects SOC concentration with DOC concentration in headwater streams, where a single parameter, the transformation rate from SOC in the soil to DOC leaching flux,  $P_r$ , accounts for the overall processes governing SOC conversion to DOC and leaching from soils (along with runoff) into headwater streams. We then derive a high-resolution  $P_r$  map over the contiguous U.S. (CONUS) using SOC data from two different sources: the Harmonized World Soil Database v1.2 (HWSD) and SoilGrids 2.0. Both maps are developed following the same five major steps: 1) selecting headwater-independent catchments where observed riverine DOC data are available with reasonable quality; 2) estimating catchment-average SOC for the independent catchments based on high-resolution SOC data; 3) estimating the  $P_r$  values for these catchments based on the generic formula and catchment-average SOC; 4) developing a predictive model of  $P_r$  with machine learning (ML) techniques and catchment-scale climate, hydrology, geology, and other attributes; and 5) deriving a national map of  $P_r$ , based on the ML model. For evaluation, we compare the DOC concentration derived using the  $P_r$  map and the observed DOC concentration values at another 3240 headwater gauges evaluation catchments. The resulting mean absolute scaled error and coefficient of determination are 0.73 and 0.47 for the HWSD-based model and 0.58 and 0.72 for the SoilGrids-based model, respectively, suggesting the effectiveness of the overall methodology. Efforts to constrain uncertainty and evaluate sensitivity of  $P_r$  to different factors are discussed. To illustrate the use of such a map, we derive a riverine DOC concentration reanalysis dataset for more than two million small catchments over CONUS. The two  $P_r$  maps, robustly derived and empirically validated, lay a critical cornerstone for better simulating the terrestrial carbon cycle in land surface and [earthEarth](#) system models. Our findings not only set a foundation for improving our predictive understanding of the terrestrial carbon cycle at the regional and global scales, but also hold promises for informing policy decisions related to decarbonization and climate change mitigation.

Formatted: Authors

35 Notice: This manuscript has been authored by UT-Battelle, LLC, under contract DE-AC05-00OR22725 with the US  
Department of Energy (DOE). The US government retains and the publisher, by accepting the article for publication,  
acknowledges that the US government retains a nonexclusive, paid-up, irrevocable, worldwide license to publish or reproduce  
the published form of this manuscript, or allow others to do so, for US government purposes. DOE will provide public access  
to these results of federally sponsored research in accordance with the DOE Public Access Plan ([https://www.energy.gov/doe-](https://www.energy.gov/doe-public-access-plan)  
40 [public-access-plan](https://www.energy.gov/doe-public-access-plan)).

## 1 Introduction

With the Earth's climate rapidly warming due to increasing atmospheric greenhouse gas concentrations, there is a growing  
focus on quantifying the regional and global carbon pools within the land, riverine, and oceanic systems, as well as the intricate  
interconnections among them (Duarte, 2017; (Jing et al., 2021; Teodoru et al., 2015; ~~Duarte, 2017~~)). Each year, about 2 billion  
45 metric tons (~~Pe~~) of dissolved organic carbon (DOC) are transported from land to the oceans via rivers globally, comparable to  
the amount of atmospheric CO<sub>2</sub> that deposits into the ocean (Hansell et al., 2009; Lønborg et al., 2020~~)).~~ Moreover, riverine  
DOC is vital to aquatic biogeochemistry by providing nutrients to microbial communities and influencing aquatic greenhouse  
gas emissions (Li et al., 2019~~)).~~

50 However, it remains a challenge to represent and predict riverine DOC effectively in the land biogeochemical module of Earth  
system models (~~ESMs~~), which are the primary tools for studying carbon cycles in the context of climate change. A chief  
reason behind this long-standing challenge is the complexity of terrestrial and aquatic processes and their interactions  
governing SOC transformation to DOC and transport from soils to rivers. The relevant terrestrial processes include the  
conversion of solid SOC into soil DOC, the adsorption and desorption of DOC by surrounding soils, the transport of DOC  
55 from soils into headwater streams along with runoff, and the degradation of soil DOC during this transport. These processes  
are further influenced by numerous biotic factors, such as microbial, plant, and enzymatic activities, as well as abiotic factors,  
including soil temperature, moisture, pH (Davidson and Janssens, 2006; Kaiser and Kalbitz, 2012; Kalbitz et al., 2000;  
Sinsabaugh, 2010). The relevant aquatic processes include the transportation of riverine DOC from headwater streams, the  
interception of DOC fluxes by reservoirs and lakes, the degradation of riverine DOC during transport, and the consumption of  
60 DOC by aquatic biosystems. Furthermore, each process is controlled by several environmental factors, which often exhibit  
substantial spatial heterogeneity. Models attempt to represent these complexities through parameters associated with governing  
equations. For instance, Tian et al. (2015a, b) incorporated the effects of runoff on DOC leaching with a coefficient that  
involves both surface and subsurface runoff. Surface and subsurface runoff are further affected by many environmental factors  
such as climate, soil, vegetation, and topography (Li et al., 2014; Li and Sivapalan, 2014~~)).~~

65

The complexity of relevant processes and their driving environmental factors is also evident in the diverse process descriptions in several land biogeochemical models that are pioneers in representing the suite of processes from SOC to riverine DOC, such as Dynamic Land Ecosystem Model (DLEM) (Tian et al., 2015a, b; Yao et al., 2021), the integrated catchment model for carbon (INCA-C) (Futter et al., 2007), the Joint UK Land Environment Simulator Dissolved Organic Carbon model (JULES-DOCM) (Nakhavali et al., 2018), and the TRIPLEX-hydrological routing algorithm (TRIPLEX-HYDRA) (Li et al., 2019). These models differ in the processes involved and the process descriptions, owing to the inconsistent understanding of relevant processes among the modeling community. For instance, DLEM and TRIPLEX-HYDRA both adopt CENTURY-like (~~Parton et al., 1987~~; Metherell et al., 1993; ~~Parton et al., 1987~~) formulas to estimate DOC leaching fluxes (Tian et al., 2015a, b; Yao et al., 2021; Li et al., 2019), but with notably different ways of incorporating both soil and water-related factors. For instance, TRIPLEX-HYDRA includes an empirical coefficient to account for soil absorption of SOC before its dissolution and DOC degradation in soils, which are not explicitly accounted for in DLEM. TRIPLEX-HYDRA incorporates hydrologic effects by directly using the water flow rate, whilst DLEM uses a dimensionless ratio to account for these effects. Equally important, the available observations have not been fully used for estimating or calibrating the numerous DOC-related parameters at the regional and larger scales in a spatially continuous yet variable fashion. Existing models usually calibrate several DOC-related parameters against DOC observations at a limited number of river ~~gauges~~stations, leading to ~~the issue of~~ overparameterization, where multiple combinations of parameter values can achieve the same simulation results (Sivapalan, 2005). Moreover, the resulting parameters often poorly reflect the spatial heterogeneity of underlying processes and environmental factors due to the limited spatial coverage of DOC observations (Futter et al., 2007; Tian et al., 2015a, b; Nakhavali et al., 2018; Li et al., 2019; Liao et al., 2019; Yao et al., 2021). Overall, existing models for simulating DOC fluxes are still subject to limited transferability over poorly observed regions due to insufficient process understanding, data scarcity, and overparameterization.

One traditional strategy for improving model transferability over poorly observed regions is parameter regionalization. Generally, ~~the~~ low-dimensional relationships between a target parameter and other environmental variables are derived based on prior knowledge or regression analysis from the locations where sufficient observations are available. The relationships are then generalized and transferred to poorly-observed places (~~Doron et al., 2011; Dupas et al., 2013; Ye et al., 2014~~; Alebachew et al., 2014; Ayata et al., 2018; ~~Doron et al., 2011; Dupas et al., 2013~~; Tan et al., 2022; ~~Ye et al., 2014~~). However, such a strategy will not work well if statistically robust and mechanistically meaningful relationships can not be derived from the conventional regression analyses or prior knowledge when, for example, the relationships are high-dimensional and nonlinear (Abeshu et al., 2022; Li et al., 2022). Fortunately, state-of-the-art machine learning (ML) techniques offer a promising and effective alternative strategy, owing to their proven advantages in capturing higher-order relationships between the target and predictive variables (~~predictors~~), especially when prior knowledge of such relationships is still in its infancy (Afan et al., 2016). For example, ML techniques have been successfully employed to capture the complex relationships ~~among~~between median sediment particle size (~~D50~~) and several environmental factors, which enabled the derivation of a national map of

median sediment particle size  $D_{50}$  (Abeshu et al., 2022). They have also been used to predict the concentration of fecal indicator bacteria, providing valuable guidance to beach closure problems (Li et al., 2022).

As the first step in addressing these challenges, this study develops an ML-powered approach for parameterizing DOC leaching fluxes at regional and continental scales. The rest of this paper is organized as follows. Section 2 outlines the overall methodology, including governing equations and corresponding parameters, data preparation, and the ML techniques employed. Section 3 presents the results over the contiguous United States (CONUS). Sections 4, 5, and 6 discuss the uncertainty, potential use of the resulting datasets, limitations of methods, and data availability. Section 7 concludes with a summary and potential future directions.

## 2 Methods

The methodology here is described with specific details over the CONUS region, but it is transferable to other regions after some modifications based on data availability.

### 2.1 Governing Equation

Several existing land or land biogeochemical models commonly employ CENTURY-like formulas to represent the leaching of DOC (Futter et al., 2007; Tian et al., 2015a, b; Nakhavali et al., 2018; Li et al., 2019; Yao et al., 2021; Parton et al., 1998). In such formulas, the DOC leaching flux is estimated as a linear function of several factors, including the SOC or DOC concentration in soil, runoff, and other relevant environmental factors. For example, in DLEM (Tian et al., 2015a, b), DOC leaching flux is estimated as

$$F_{DOC\_runoff} = F_{SOC\_soil} \times \alpha_1 \times \alpha_2 \times \alpha_3 \quad (1)$$

Where  $F_{SOC\_soil}$  is the total amount of decomposed SOC in soil ( $\text{g Cm}^{-2}\text{s}^{-1}$ );  $\alpha_1$  is the fraction of decomposed SOC that is dissolvable (%);  $\alpha_2$  is the runoff coefficient (-), i.e., the ratio of total runoff volume to the sum of total runoff volume and soil water content; and  $\alpha_3$  is another coefficient (-) accounting for the effects of DOC concentration in soil water and desorption. In TRIPLEX-HYDRA (Li et al., 2019), DOC leaching flux is given as

$$F_{DOC\_runoff} = C_{SOC} \times K_s \times K_a \times Q_{runoff} - K_{soil} \quad (2)$$

where  $F_{DOC\_runoff}$  is the DOC flux in the soil water ( $\text{g C/s}$ );  $C_{SOC}$  is the concentration of SOC in the soil ( $\text{g C/m}^3$ );  $K_s$  is the solubility of SOC (-);  $K_a$  is the adsorption coefficient of SOC (-);  $K_{soil}$  represents the degradation rate of DOC in soils ( $\text{g C/s}$ ), and  $Q_{runoff}$  is total runoff rate ( $\text{m}^3/\text{s}$ ).

Based on the similarity between equations (1) and (2), while keeping minimal complexity in the process representation, we propose a simpler formula to estimate DOC leaching flux as

$$F_{DOC\_runoff} = C_{SOC} \times Q_{runoff} \times P_r \quad (3)$$

Eqn. (3) can be rewritten as

$$C_{DOC\_runoff} = \frac{F_{DOC\_runoff}}{Q_{runoff}} = C_{SOC} \times P_r \quad (4)$$

Where  $F_{DOC\_runoff}$  is the DOC leaching flux (g C/s),  $C_{SOC}$  is the SOC concentration (g C/m<sup>3</sup> soil),  $Q_{runoff}$  is the runoff volume per unit time (m<sup>3</sup> water/s),  $P_r$  is the transformation rate from SOC in soil to DOC in runoff (m<sup>3</sup> soil/ m<sup>3</sup> water), and  $C_{DOC\_runoff}$  is the DOC concentration in the runoff (g C/m<sup>3</sup> water).

Eqn. (4) has ~~several~~two advantages: 1) its lumped parameter,  $P_r$ , accounts for all relevant processes and factors, including soil carbon decomposition, DOC sorption-desorption balance, DOC transport and degradation in soils, etc.; 2) its simplicity significantly reduces data requirements for large-scale parameterization since it is highly parameter-parsimonious and much more compatible with the availability of DOC observational data.

For a "small catchment", we further assume that  $C_{DOC\_runoff}$  can be approximated with the riverine DOC concentration at the catchment-outlets for headwater catchments, i.e.

$$C_{DOC\_outlet} \approx C_{DOC\_runoff} \quad (5)$$

Where  $C_{DOC\_outlet}$  is the riverine DOC concentration at the catchment outlet (g C/m<sup>3</sup>). In this study, a "small catchment" refers to the drainage basin extending from the river station upstream to the furthest tributaries that do not have any upstream rivers. Note that a small catchment is not necessarily a headwater catchment that includes only one river (He et al., 2024). The rationale behind Eqn. (5) is two-fold: 1) The travel time of runoff in small headwater-streams of small catchments is typically much less than one day, e.g., the daily total runoff rate can be approximated with the daily streamflow rate for headwater small catchments (Li et al., 2013; Ducharme et al., 2003; Li et al., 2013), and 2) Due to the short travel time of DOC in headwater streams, riverine DOC degradation in headwater streams mostly occurs at a rate of about 1% per day according to previous, based on our literature review of existing experimental (Qualls and Haines, 1992; Sobczak et al., 2003) and modeling studies (Strauss & Lamberti, 2002; Tian et al., 2015a, b; Li et al., 2019; hence is) (for a full list of references, see Supplementary Table S1). Given this minimal degradation rate and the short residence time of DOC in streams of small catchments (on the order of a few hours), it is reasonable to assume negligible stream to the point it exits into downstream rivers. Combining Eqn. (4) and (5) yields

$$C_{DOC\_outlet} \approx C_{SOC} \times P_r \quad (6)$$

Eqn. (6) may be used in at least two ways: 1) One can estimate  $P_r$  at the catchment scale wherever observed DOC concentration and SOC values are available, and 2) Once  $P_r$  is estimated a priori or through calibration, one can quickly predict riverine DOC concentration or discharge in headwater-streams of small catchments from the corresponding SOC values.

## 2.2 Data

A key step in the data preparation in this study is to pair up SOC data and riverine DOC observations at headwater catchments. The SOC data required for this study are from the Harmonized World Soil Database (HWSD) v1.2 (Fischer et al., 2008). This database provides SOC values at a spatial resolution of 1 km for two vertical soil layers at each grid cell—the top layer (0–30 cm) and the sub-layer (30–100 cm). Considering that DOC leaching from soils into rivers predominantly comes from the topsoil (Brooks et al., 1999; Finlay et al., 2006), we use the SOC content data from the top 30 cm layer for our estimations. We also take into consideration that there are missing values in some grid cells in the HWSD v1.2 and adjust our catchment selection accordingly. Riverine DOC observations are available via the Water Quality Portal (WQP) (Water Quality Portal, 2021). WQP integrates the publicly available water quality data from the USGS National Water Information System (NWIS) (U.S. Geological Survey), the EPA STORage and RETrieval Water Quality eXchange (STORET-WQX) (USEPA), and the USDA ARS Sustaining The Earth's Watersheds - Agricultural Research Database System (STEWARDS) (Steiner et al., 2008). As of now, the WQP features data from 32071 river stations within the CONUS. These stations have recorded at least one DOC measurement between 1900 and the present.

Regional and global soil property maps, such as soil organic carbon (SOC) maps, are typically generated using two primary methods: the linkage method (also known as the taxotransfer rule-based method) (Batjes, 2003) and digital soil mapping (McBratney et al., 2003). This study employs the most widely recognized datasets from each method: the Harmonized World Soil Database (HWSD) v1.2 (Fischer et al., 2008) and SoilGrids 2.0 (Poggio et al., 2021). HWSD provides SOC data at a spatial resolution of 1 km for two soil layers—the top layer (0–30 cm) and the sub-layer (30–100 cm). As one of the first globally harmonized soil datasets, it integrates data from diverse national and regional sources into a standardized framework, making it a foundational resource for many Earth system modeling studies (Best et al., 2011; Han et al., 2014; Todd-Brown et al., 2013; Zhao et al., 2018). SoilGrids 2.0 offers SOC data at a higher resolution of 250 m for the same layers, leveraging machine learning algorithms to enhance accuracy and constrain uncertainty. Its higher resolution and improved reliability have made it increasingly popular for Earth system modeling since its release (Dai et al., 2019; Hengl et al., 2017; Poggio et al., 2021). Considering that DOC leaching from soils into rivers predominantly comes from the topsoil (Brooks et al., 1999; Finlay et al., 2006), we use the SOC content data from the top 30 cm layer for our estimations. We also take into consideration that there are missing values in some grid cells in the HWSD v1.2 and SoilGrids 2.0 and adjust our catchment selection accordingly.

In order to pair up SOC and DOC data at headwatersmall catchments, we rely on the National Hydrography Dataset Plus (NHDPlus) dataset hosted by the U.S. Geological Survey (USGS) (Mckay et al., 2012). This dataset is chosen for two reasons: Firstly, NHDPlus provides well-defined catchment boundaries and their-correspondingassociated river segments, denotedreferred to as local catchments and flowlines. There-areIt includes ~2.6 million NHDPlus-flowlines in across CONUS,

drainage catchment for any flowline, which is the sum of both local catchment and the drainage areas corresponding to all the flowlines upstream of the local one, can be derived from the established flowline network. The sizes of these 2.6 million local catchments vary from the 5<sup>th</sup> percentile at 9.680.02 km<sup>2</sup> to the 95<sup>th</sup> percentile at 0.029.68 km<sup>2</sup>, depending on the corresponding surface topography, with a CONUS average of 3.12 km<sup>2</sup> (see supplementary Figure Supplementary Fig. S1). Secondly, NHDPlus is closely linked to ScienceBase (Wieczorek et al., 2018), a comprehensive scientific data and information management platform also hosted by USGS. ScienceBase incorporates includes a wide range of environmental variables, including across 11 categories, such as climate, hydrology, soil, and geological data, conveniently available at the catchment scale over across the whole entire CONUS. These environmental data are critical in the ML modeling analysis.

Correspondingly, the overall data preparation procedure consists of three major steps: 1) Selection of headwaters small catchments based on the availability of observed riverine DOC concentrations of adequate quality. 2) Estimation of  $P_r$  values for the catchments selected in Step 1, leveraging the corresponding riverine DOC observations and SOC reanalysis data. 3) Extraction of catchment-scale environmental variables that could potentially influence  $P_r$ . Specific details of each step will be further discussed in the following subsections. This study adopts two SOC datasets, both of which directly influence the calculated  $P_r$  values used in training, thereby affecting all steps leading to the final  $P_r$  map. To enhance clarity and avoid redundancy, the HWSO-based model is the primary focus of discussion, as the workflow and major conclusions remain consistent. More information on the SoilGrids-based model is available in the supplementary materials. Users can choose their preferred  $P_r$  map based on their specific needs.

### 2.2.1 Selecting headwaters small catchments

Our selection process for suitable headwaters small catchments involves the integration of the NHDPlus dataset and observed riverine DOC concentration data from river stations:

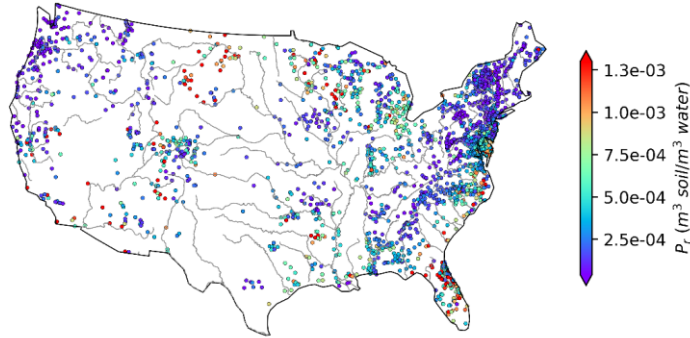
1. We conduct a geospatial analysis to identify the upstream drainage area of each WQP river station. This is accomplished by using the NHDPlus local catchments and flowlines. For every WQP station, we search for a NHDPlus flowline which the station is located. Using the Python package HyRiver (Chen et al., 2021), we locate 29,320 located 29,320 WQP stations with the closest corresponding NHDPlus flowlines. However, the remaining 2751,2751 stations cannot be linked with the NHDPlus dataset due to the absence of adjacent flowlines. Some When WQP stations are in close proximity to each other and share the same NHDPlus flowlines. In such a case flowline, we retain only one WQP the station with the best data availability. Each For a given flowline in NHDPlus is accompanied by a corresponding watershed boundary. However, not all WQP stations are precisely located at the outlets of these existing NHDPlus watershed, HyRiver traces it back to every upstream flowline, accessing and merging the boundaries. When faced with these circumstances, we derive the upstream drainage area of all related NHDPlus local catchments from the Hydro Network-Linked Data Index web server. It also requests the server to simplify the boundaries for the WQP stations from Digital Elevation Model (DEM) data. Upon completion

of and split them precisely at the station locations. The relationship between the derived small catchment boundaries and the NHDPlus local catchments is shown in Supplementary Fig. S2a. Through this comprehensive geospatial analysis, we identify the upstream boundaries for 22,201 WQP stations.

2. We further select the WQP stations whose drainage areas can be considered headwaterssmall catchments, based on two criteria: 1) there are no upstream rivers flowing into them, and 2) their drainage areas are no more than 2500 km<sup>2</sup>. This size threshold ensures that the travel distance of river water (and consequently, DOC) is ~50 km within these catchments. Assuming an average channel velocity of ~1.0 m/s (Chow et al., 1988), the average travel time is ~14 hours, i.e., less than one day. Using these criteria, we identify 18,612 pairs of WQP stations and headwaterssmall catchments.
3. For the 18,612 WQP stations, we perform a rigorous DOC data quality control based on five criteria: a) The record lengths of riverine DOC data should span at least one year; b) There should be at least two riverine DOC observations; c) No single season should dominate the riverine DOC observations, i.e., a single season should not account for more than 50% of the records; d) within the boundaries of the corresponding catchments, there should be sufficient availability of the NHDPlus catchment attributes and SOC reanalysis data; e) the catchments should not be significantly affected by dams, i.e., the total drainage areas of the dams within a catchment should be no more than 5% of the total catchment area. The adoption of criteria (a)-(e) reflects a careful balance between ensuring data quality and maintaining adequate quantity, ensuring that sufficient WQP stations are retained to represent the entire CONUS. After the data quality control, there remain 5805 WQP stations with their corresponding headwaterssmall catchments.
4. For the 5805 WQP stations and their headwaterssmall catchments, we verify the spatial independence among them. ~~For instance, Catchment A-A catchment is considered to be nested within Catchment B-another if A is situated it lies entirely within the latter's drainage area of B. In such scenarios, while. While the fluxes-observed flux at the downstream catchment's outlet of Catchment B are dependent depends on those at the outlet of Catchment A, Catchment A itself remains independent of B-contributions from upstream catchments, the upstream catchments maintain hydrological independence. As stated in Supplementary Figure 2, here of catchment headwaters, the same logic as Fig. 3a in nesting scenario shows two gray catchments, A and B, both located within the red catchment, C. Since A and B have no containing relationship and are both smaller than C, they are classified as independent catchments. In contrast, C is considered a nesting catchment. The same logic applies consistently selected as the independent catchment in more complex nesting scenarios. From the 5805 pairs of the WQP stations and catchments, we identify 2595 as being independent and suitable for further ML modeling-model training. The other 3210 pairs, despite the nesting issue, are still valuable; they are thus kept for evaluation of estimated DOC (see Section Sect. 3.4). Due to missing values in SoilGrids 2.0, valid  $P_r$  estimates are unavailable for 12 out of 2595 independent catchments; however, the number of evaluation catchments remains unchanged.~~



a)  $P_r$  of independent catchments



b)  $P_r$  of evaluation catchments

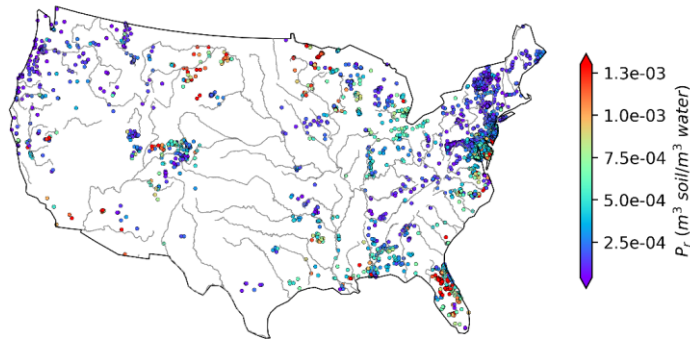


Figure 1. Variability in estimated  $P_r$  across CONUS: a) For independent catchments ( $n=2595$ ), and b) For evaluation catchments ( $n=3210$ ). The points indicate the locations of the WQP stations, which are also the outlets of the corresponding small catchments. The CONUS boundary and river shapefiles are directly obtained from open-source datasets GeoPandas (geopandas.org) and Natural Earth (Made with Natural Earth. Free vector and raster map data @ [naturalearthdata.com](https://www.naturalearthdata.com)), respectively. The color bars have been adjusted to enhance visual display by showing only the main body of values (from the 5th percentile to the 95th percentile).

### 2.2.2 Estimating $P_r$

For the final set of the paired WQP stations and headwaterssmall catchments, we calculate  $P_r$  using the DOC observation from the WQP stations and long-term mean SOC from HWSD based on Eqn. (6). For each catchment, the catchment polygons are

used to clip the top-layer SOC map at the 1km resolution, and the catchment-scale SOC is subsequently calculated as the spatial average of SOC values at those 1km grid cells within the catchment. Hereafter the  $P_r$  estimated using Eqn. (6) are referred to as "Estimated  $P_r$ ". The Estimated  $P_r$ , derived from the analysis of WQP DOC observations and HWSO SOC data, exhibits a wide range of values spanning several orders of magnitude. Figure 1a illustrates the spatial distribution of  $P_r$  for the 2595 independent catchments. In these catchments, the Estimated  $P_r$  ranges from  $4.61 \times 10^{-6}$  to  $8.04 \times 10^{-3}$  ( $\text{m}^3 \text{ soil} / \text{m}^3 \text{ water}$ ), with a median value of  $2.50 \times 10^{-4}$  ( $\text{m}^3 \text{ soil} / \text{m}^3 \text{ water}$ ). As a broad assessment of the similarity between the catchments used to construct the model and the evaluation catchments, the values of  $P_r$  for the evaluation catchments calculated from data values of DOC and SOC using Eqn. (6) are shown in Figure 1b. Here, the Estimated  $P_r$  values in these catchments range from  $8.81 \times 10^{-6}$  to  $6.37 \times 10^{-3}$  ( $\text{m}^3 \text{ soil} / \text{m}^3 \text{ water}$ ), with a median of  $2.60 \times 10^{-4}$  ( $\text{m}^3 \text{ soil} / \text{m}^3 \text{ water}$ ). Note that the spatial distribution of the selected catchments is quite consistent with the spatial distribution of the WQP stations, i.e., more densely distributed in the eastern than western U.S., suggesting a good spatial representation of the selected catchments over all the WQP stations in CONUS. Figure S8 shows the spatial distribution of Estimated  $P_r$  values derived from the SoilGrids-based model for independent and evaluation catchments. The overall pattern closely resembles that derived from the HWSO-based model. The Estimated  $P_r$  values have a slightly narrower range, from  $1.16 \times 10^{-5}$  to  $8.69 \times 10^{-3}$  ( $\text{m}^3 \text{ soil} / \text{m}^3 \text{ water}$ ) at independent catchments, and a similar range, from  $7.78 \times 10^{-6}$  to  $7.55 \times 10^{-3}$  ( $\text{m}^3 \text{ soil} / \text{m}^3 \text{ water}$ ) at evaluation catchments.

### 2.2.3 Extracting environmental variables

The ScienceBase dataset is a comprehensive resource that houses a wide array of environmental variables sorted into categories such as climate, hydrology, geology, and land use/land cover. We collect a wide range of environmental variables, comprising a total of 126 variables, across eleven from the ScienceBase dataset, spanning 11 distinct categories. We remove seven attributes related to dams and streams from the analysis as they are excluded as irrelevant to our analysis objectives. Furthermore, we exclude, along with 24 attributes from further analysis because they contain predominantly zero values, with over (>80% of the values being zero over %) across CONUS. Out of the remaining 95 variables (see supplementary Tables S1 and S2 for details), 46 are relatively independent from each other. However, the other while 49 are highly correlated showed strong correlations with one or more variables. These 49 non-independent variables are further Following Schober et al. (2018), we define strong correlation as a Pearson correlation coefficient  $|r| \geq 0.8$ . The 49 correlated variables are categorized into 9 distinct "correlated groups" and named based on the group property, as listed in Table 1. A "correlated group" is characterized by shared properties, where each variable demonstrates a strong correlation with at least one other variable within its group but a weak correlation ( $|r| < 0.8$ ) with variables outside the group. We address the interdependence within each "correlated group" in through two steps. First, we normalize each variable within a group: 1) normalizing individual variables using the Yeo-Johnson power transformation (Yeo and Johnson, 2000) (see Supplementary Figure S3). The transformation ensures that the resulting dataset has a to achieve zero mean of 0.0 and a unit variance of 1.0. Second, we merge all (Supplementary Fig. S3), and 2) merging the normalized variables into a single new variable through linear summation to create a single new variable (Daoud, 2018). This new variable is thusnow relatively independent of the

other environmental variables. For those 46 ~~independent~~ variables, we apply the same transformation to minimize the impacts of varying magnitudes between different variables. Eventually, 54 ~~independent~~ variables remain, including 46 originally relatively independent and 9 newly merged variables from the correlation groups—(see Supplementary Tables S2 and S3 for details).

~~The ML technique used in this study is~~We use the eXtreme Gradient Boosting (XGBoost) ~~algorithm,~~, which is a powerful and widely adopted ~~machine learning~~ML algorithm due to its exceptional performance in various applications (Abeshu et al., 2022; Delavar et al., 2019; Li et al., 2022). XGBoost is a scalable end-to-end tree-boosting system that belongs to the ensemble learning family-(Chen and Guestrin, 2016). It combines multiple weak learners into a strong learner via sequential training and improving, and eventually forms a robust and accurate predictive model. By using XGBoost in this study, we aim to develop a predictive model that establishes causal linkages between the target variable,  $P_r$ , and a small number of environmental variables (denoted as predictors hereafter).

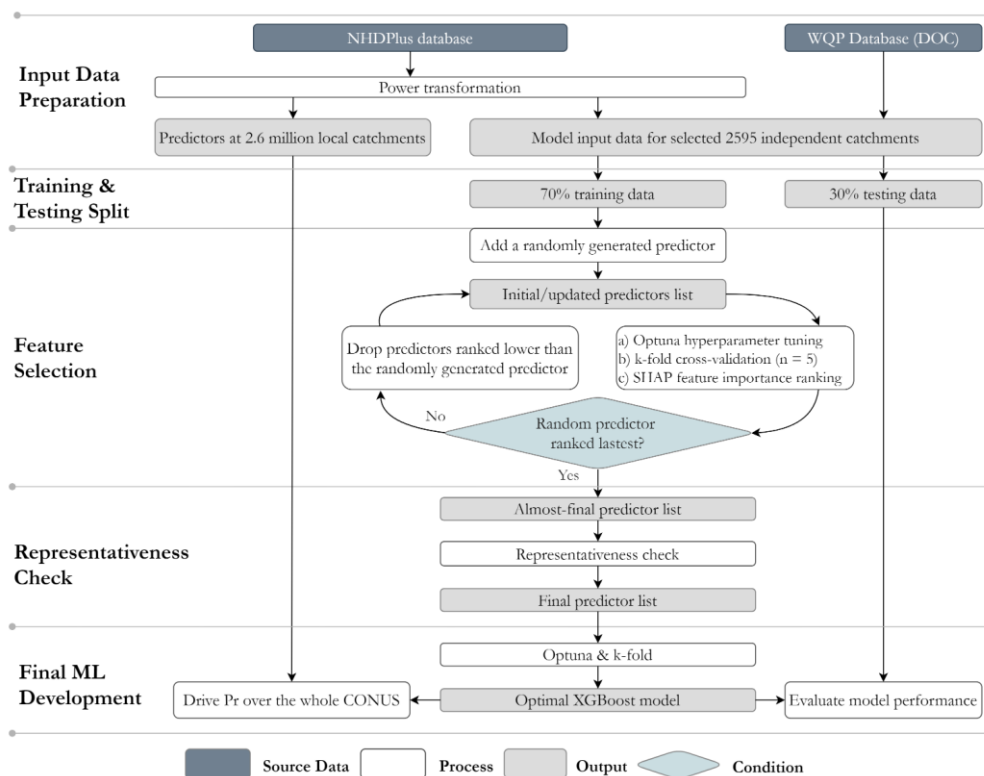
In addition to XGBoost, we take advantage of some other ML tools and techniques. Specifically, we use the Optuna optimization framework (Akiba et al., 2019) and k-fold cross-validation (k=5) for tuning the hyperparameters. By leveraging Optuna and k-fold cross-validation, we can systematically search and optimize the hyperparameters, maximizing the model's performance and accuracy. Furthermore, we employ the SHapley Additive exPlanations (SHAP) (Lundberg and Lee, 2017) to aid in the selection of environmental factors that are related to  $P_r$ . SHAP is a technique that assigns importance values to individual predictors in a model, providing insights into their contributions to the prediction. By using SHAP, we can identify the key environmental factors that significantly influence  $P_r$  and further refine our model. ~~Recent studies have demonstrated the efficiency and effectiveness of these techniques in capturing high-dimensional and complex relationships between a target biogeochemical variable and various environmental predictors.~~These techniques have been successfully applied in various studies, including riverine sediment, beach water quality, oceanic particulate organic carbon, and eutrophication impacts from corn production (Abeshu et al., 2022; Fan et al., 2021; Li et al., 2022; Liu et al., 2021; Romeiko et al., 2020; Fan et al., 2021), demonstrating their efficiency and effectiveness in capturing high-dimensional and complex relationships between a target biogeochemical variable and various environmental predictors. Readers are referred to Abeshu et al., (2022) for more details about these techniques.

The overall procedure for developing a predictive ML model is illustrated in ~~Figure~~Fig. 2 (identical for the SoilGrids-based model) and outlined as follows:

1. Prepare the input data for the ML modelling based on the independent catchments, their corresponding  $P_r$  estimates, and environmental variables. To address the substantial statistical disparities and wide variation within each predictor, we employ power transformation on all predictors. The lambda parameter is held constant during the transformation process for the training, testing, and prediction datasets to ensure consistent and reproducible results. Following the

transformation, the dataset exhibits a zero-mean and unit variance, with a distribution that closely resembles a Gaussian distribution (as illustrated in Figure Supplementary Fig. S3).

2. Randomly split the observational dataset (2595 catchments) into two sets: 70% for training and 30% for testing the ML model. These training and testing sets will be used throughout the subsequent steps.
3. Identify the list of predictors out of the 54 environmental variables extracted in Section Sect. 2.2.3 in three sub-steps:
  - a. Generate a completely random predictor.
  - b. Prepare an initial list of candidate predictors consisting of the random predictor and an initial list of candidate environmental variables. Use Optuna and k-fold cross-validation to obtain the optimal hyperparameters and train an intermediate ML model until the model achieves the best performance evaluated using the testing set.
  - c. Calculate and rank the SHAP values for all the candidate predictors. Update the list of candidate predictors by keeping only those predictors with better SHAP values than the random predictor. For example, if the random predictor is ranked 20th, only the top 19 predictors are passed to the next iteration.
  - d. Obtain an almost-final list of predictors by repeating sub-steps b-c.
4. Check the representativeness of the almost-final list of predictors identified in Step 3. For each of these predictors, check whether its values from the independent catchments are statistically representative of the whole CONUS, i.e., its values from those 2.6 million local catchments. Drop those predictors that cannot pass the representativeness check. Similar to Abeshu et al. (2022), the representativeness check on each of the almost-final predictors is performed by comparing the cumulative distribution function (CDF) derived from the observational dataset (2595 training catchments) and the CDF derived from the whole CONUS (about 2.6 million local catchments in NHDPlus). Specifically, comparisons are made between the 5th, 25th, 50th, 75th, and 95th percentiles between the two CDFs. After this Step 4, a final list of predictors is obtained.
5. Develop the final ML model based on the final list of predictors using Optuna and k-fold cross-validation methods.



**Figure 2. A workflow for the XGBoost model.**

In Steps 3 and 5, model performance metrics are required for model training and validation/evaluation. The Kling-Gupta efficiency (KGE) (Gupta et al., 2009) has the advantage of simultaneously capturing both the magnitude and phase differences between the observed and simulated series ((Gupta et al., 2009; Abeshu et al., 2022); Gupta et al., 2009). However, further investigations have revealed several limitations: a) lack of an inherent benchmark value to distinguish between "good" and "bad" model performance, b) sensitivity to outliers, which can result in a systematic overestimation of the target variable, and c) instability when the target variable approaches zero (Knoben et al., 2019; Pool et al., 2018; Santos et al., 2018; Knoben et al., 2019). Therefore, in addition to KGE, the mean absolute scaled error (MASE) is also used here to alleviate the influence of extreme values in the observation or simulation data (Hyndman and Koehler, 2006). MASE is a scaled error metric that is defined as the mean absolute error (MAE) of the model simulation divided by scaling factors (MAE of the observation in the

original definition). In this study, we normalize MAE by the geometric mean of the observation data. Note that Steps 3 and 5 above are relatively independent of each other and do not have to rely on the same metrics.

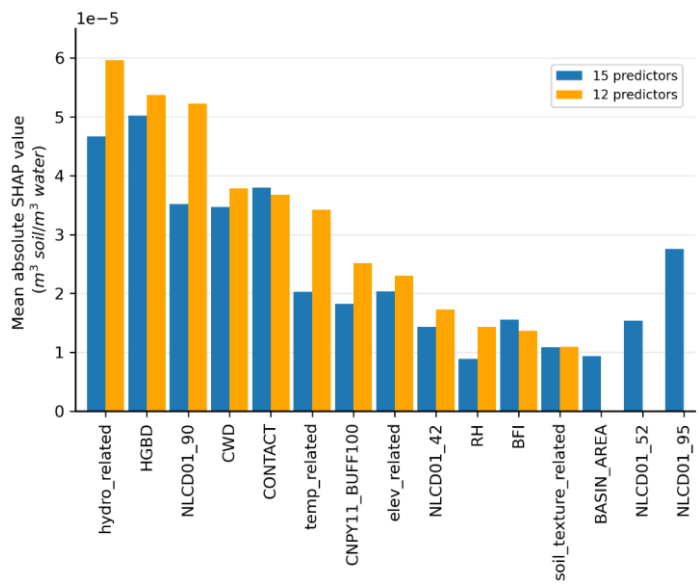
3 Results

3.1 Predictor selection

In the predictor selection stage, after six iterations of hyperparameter tuning and predictor reduction with KGE as the metric, a list of 15 predictors is selected (see Table 2 blue bars in Fig. 3), including those related to climate, hydrology, pedology, and land cover. In addition, using MASE as the metric in this stage leads to a list of 19 remaining predictors, among which 13 are the same as the list of predictors identified using KGE. The predictor list selected using KGE is preferred due to the fewer predictors and similar model performance. Figure S9 shows the feature selection results (blue bars) for the SoilGrids-based model, with 11 out of 13 predictors also included in the final list derived from the HWSO-based model. This overlap further reinforces the consistency of important features across datasets and enhance the robustness of the selection process.

To further investigate the model performance, we conducted a sensitivity analysis using the 15 predictors. We selected 3 of 3 initially selected predictors: "BASIN\_AREA," "NLCD01\_52," and "NLCD01\_95." These variables demonstrated insufficient representativeness of the anticipated real-world data distribution in the prediction phase, resulting in a final model with 12 predictors. Figure 3 presents a comparative analysis of mean absolute SHAP values between the original 15-predictor model (blue bars) and the final 12-predictor model (orange bars). Notably, both models identified the same five dominant predictors, ranked according to their influence in the 12-predictor model: 1) the merged predictor of hydrologic variables ("hydro\_related"), 2) the areal percentage of Hydrologic Group BD soil ("HGBD"; detailed classification in Ross et al., 2018), 3) the areal percentage of woody wetlands ("NLCD01\_90"), 4) the consecutive wet days ("CWD"), and 5) the subsurface flow contact time ("CONTACT"). The "hydro\_related" and "CWD" reflect the overall hydrology condition of a catchment, including runoff, precipitation, and groundwater recharge. Groundwater has a dilution effect on DOC concentration (Kortelainen and Karhu, 2006). Similarly, precipitation and runoff contribute to the distribution and concentration of DOC in rivers (Autio et al., 2016; Camino-Serrano et al., 2014; Autio et al., 2016). Woody wetland, as one land cover attribute, has been identified as a significant predictor of downstream DOC concentration (Duan et al., 2017), because of the enhanced breakdown of organic matter in the soil. The subsurface flow contact time ("CONTACT") is a key factor in the DOC concentration, as it influences the residence time of DOC in the subsurface, affecting its degradation and transformation. For instance, during the prediction phase. Therefore, only 12 predictors are adopted in the final model training-transport, a catchment with a shorter contact time experiences reduced mineralization loss (Ludwig et al., 1996) and microbial consumption (Helton et al., 2015). Conversely, studies have shown that labile DOC concentration increases with contact time in some

alluvial aquifers, as deeper groundwater inflow could provide considerable labile DOC (Helton et al., 2015; Wickland et al., 2012).

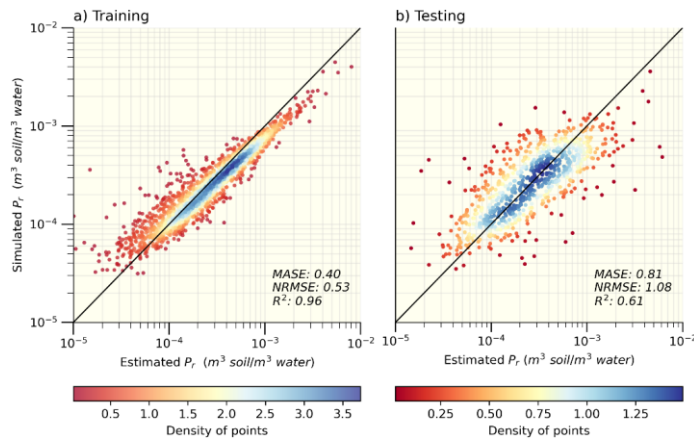


**Figure 3. Mean absolute SHAP values of predictors in models with 15 predictors (blue) and 12 predictors (orange).** Note that the SHAP values have the same units as the target variable,  $P_r$ . Abbreviations: hydro\_related (merged predictor representing recharge, runoff, and precipitation); HGBD (areal percentage of Hydrologic Group BD soil); NLCD01\_90 (areal percentage of woody wetlands); CWD (consecutive wet days); CONTACT (subsurface contact time); temp\_related (merged predictor encompassing potential evapotranspiration, first/last freeze timing, snow fraction, actual evapotranspiration, and mean/min/max temperature); CNPY11\_BUFF100 (areal percentage of canopy in the riparian buffer); elev\_related (merged predictor for mean/min/max elevation); NLCD01\_42 (areal percentage of evergreen forest); RH (relative humidity); BFI (base flow index); soil\_texture\_related (merged predictor for silt and sand content); BASIN\_AREA (catchment area); NLCD01\_52 (areal percentage of shrub); NLCD01\_95 (areal percentage of herbaceous wetlands). For detailed descriptions, refer to Supplementary Tables S2 and S3.

### 3.2 Final model

Figure 34 presents the performance of the ML model during both the training and testing phases (phases shown in FigureFig. 2). To mitigate over-plotting, all the scatter plots (Figure-3Fig. 4 and hereinafter) employ color coding based on estimated density using kernel density estimation (KDE), as indicated by the corresponding color bar. After the exclusion of the three variables that displayed poor representativeness, the ML model performance remains stable between the training and testing phases, as gauged by metrics such as MASE, coefficient of determination ( $R^2$ ), and normalized root-mean-square-error

(NRMSE). The similarities in these metrics between the *Estimated* and predicted  $P_r$  values across both phases support the robustness of our 12-predictor model. Consequently, the final ML model and the subsequent analyses are based on the 12 selected predictors. Furthermore, the consistency of model performance between the training (MASE= 0.40) and testing (MASE= 0.81) phases suggests that the model overfitting issues are well-regulated (Ying, 2019). We also use KGE as the metric during the final model training. After a comparison between the modeling results using MASE (Figure 3 Fig. 4) and KGE (supplementary Figure Supplementary Fig. S4), MASE is preferred for two reasons: a) using MASE yields a better consistency in model performance between the training and testing phases, suggesting better model transferability; b) using MASE leads to a closer agreement between the model simulated and *Estimated*  $P_r$  values. Figure S10 illustrates the performance of the SoilGrids-based model, showing similar metrics overall. However, during the testing phase (Supplementary Fig. S10b), the model slightly overestimates low values and underestimates high values. This discrepancy is likely due to the flatter data distribution in the testing dataset, which results in insufficient learning for those extreme values.



**Figure 34.** Performance of the XGBoost model with 12 predictors during a) the training phase (n=1816) and b) the testing phase (n=779). The solid black line indicates a 1:1 ratio. The varying colours indicate the density of points in the scatter plot.

Table 31 lists the optimized hyperparameter values of the final XGBoost model. (Supplementary Table S4 for that of SoilGrids-based model). We choose to tune 8 model parameters, which are critical to the XGBoost tree booster controlling regularization, subsampling, learning process, and the growth of the tree. The optimal values of model hyperparameters are quite different from the default ones, suggesting hyperparameter tuning is necessary.

Table 31. The optimal values of the XGBoost model hyperparameters.



Hyperparameter	Optimal Value	Tuning Range	Default value	Description
lambda	$6.725 \times 10^{-1}$	[0, $\infty$ ]	1	Control L1 and L2 regularization; the larger the value, the more conservative the model will be
alpha	$7.484 \times 10^{-2}$	[0, $\infty$ ]	0	
gamma	$1.316 \times 10^{-2}$	[0, $\infty$ ]	0	Govern the model learning process by changing the step size shrinkage and minimum loss reduction; the larger the value, the more conservative the model will be
eta	$1.277 \times 10^{-1}$	(0, 1]	0.3	
colsample_bytree	$9.323 \times 10^{-1}$	(0, 1]	1	Control the subsample ratio of columns and training instances; a proper set of those values will prevent the model from over-fitting
subsample	$6.142 \times 10^{-1}$	(0, 1]	1	
min_child_weight	$8.410 \times 10^{-2}$	[0, $\infty$ ]	1	Determine the growth of the tree
Maxmax_depth	12	[0, $\infty$ ]	6	

Formatted Table

Figure 45 depicts the correlation between  $P_r$  and the 12 predictors and among the predictors themselves, (Supplementary Fig. S11 for that of SoilGrids-based model), where highly positive correlated and negative correlated are shown in dark-red and blue colors, respectively. Since we have treated the highly correlated variables, the highest positive correlation coefficient is 0.63 between "per\_canopyCNPY11\_BUFF100" and "hydro\_related", lower than the threshold of 0.8 we adopt in Sect. 2.2.3. Among the observed correlation coefficients, the highest negative correlation coefficient, -0.69, is found between the variables "elev\_related" and "temp\_related." This strong negative correlation makes intuitive sense since air temperature decreases with increasing elevation. Note that all of the 12 selected predictors show weak or even negligible correlation with the target variable  $P_r$ , with the absolute values of the correlation coefficient less than 0.3. It is not surprising since the high-order, nonlinear relations among  $P_r$  and the predictors, and likely among the predictors themselves, can only be effectively captured by the ML techniques but not the traditional regression analysis methods.

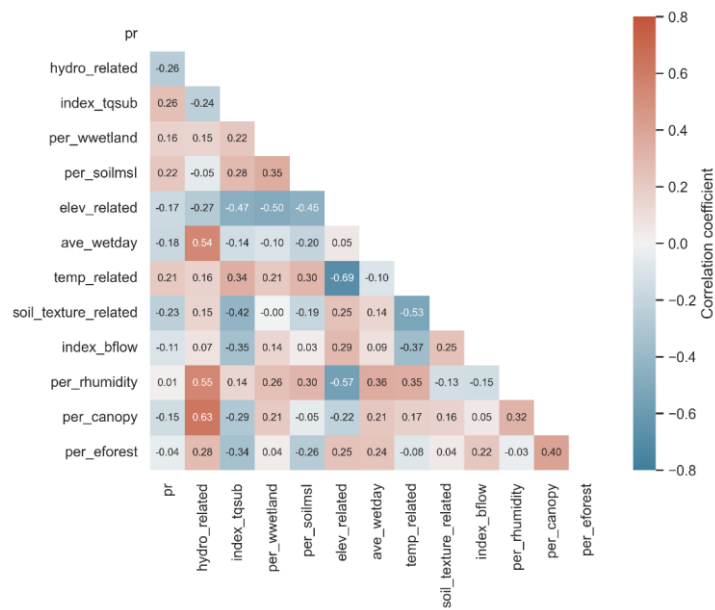


Figure 5. Covariance heatmap of  $P_r$  and the 12 selected NHDPlus predictors. The Pearson correlation coefficient is used. Abbreviations: hydro\_related (merged predictor representing recharge, runoff, and precipitation); CONTACT (subsurface contact time); NLCD01\_90 (areal percentage of woody wetlands); HGBD (areal percentage of Hydrologic Group BD soil); elev\_related (merged predictor for mean/min/max elevation); CWD (consecutive wet days); temp\_related (merged predictor encompassing potential evapotranspiration, first/last freeze timing, snow fraction, actual evapotranspiration, and mean/min/max temperature); soil\_texture\_related (merged predictor for silt and sand content); BFI (base flow index); RH (relative humidity); CNPY11\_BUFF100 (areal percentage of canopy in the riparian buffer); NLCD01\_42 (areal percentage of evergreen forest). For detailed descriptions, refer to Supplementary Tables S2 and S3.

### 3.3 $P_r$ map

We develop a spatially continuous map of  $P_r$  over CONUS by applying the final XGBoost model over the 2.6 million NHDPlus local catchments, as shown in Figure 5Fig. 6. The spatial patterns of  $P_r$  are generally consistent with those in FigureFig. 1. High  $P_r$  values, shown in orange and red, are mostly located on the southeast coasts, New Mexico, Arizona, southern California, and North Dakota. Low  $P_r$  values, shown in blue and purple, are more prevalent in the Northeast and Northwest regions. This consistency between FiguresFig. 1 and 5Fig. 6 again confirms that the 2595 independent catchments used in the ML modeling are representative of the whole CONUS domain, hence supporting the transferability of the ML modeling results.

Figure S12 presents the spatial maps derived using the SoilGrids-based model. The overall patterns are very similar at most

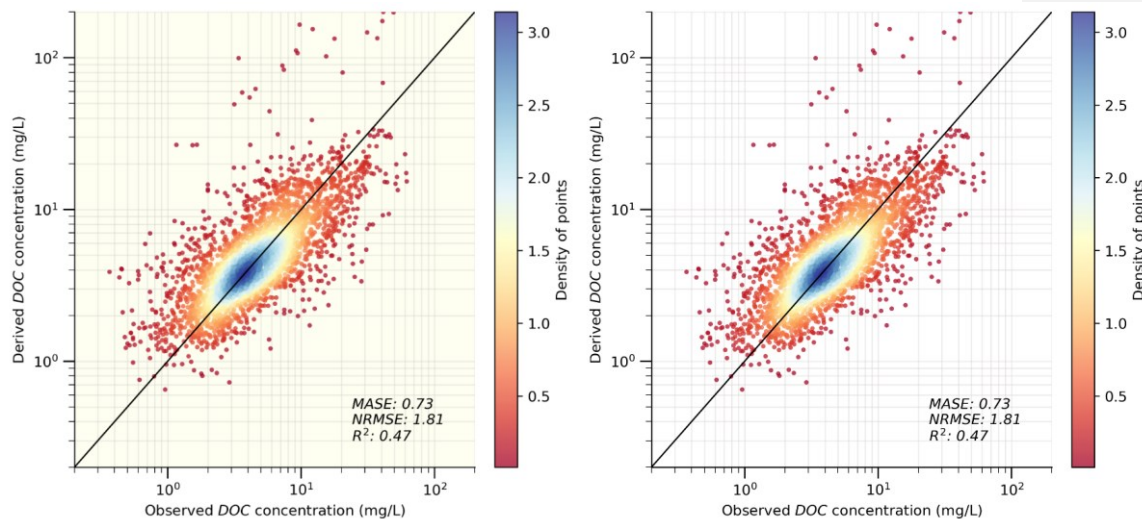
places; however, the model predicts lower values in southern California, New Mexico, and Colorado, and higher values in northern Minnesota and southern Florida.

Figure 56. ML model simulated  $P_r$  at over 2.6 million NHDPlus local catchments.

### 3.4 Evaluation

We evaluate the  $P_r$  map by comparing the DOC concentration values derived from this map (and Eqn. 6) with those observed, since there is no direct measurement of  $P_r$ . The 3210 evaluation gauges (stations and their corresponding headwater small catchments. See (Fig. 1b) are used for this purpose. Note that each of these 3210 evaluation catchments may encompass multiple NHDPlus local catchments. The evaluation thus takes two steps: 1) For each evaluation-NHDPlus local catchment, calculate its average DOC concentration using the predicted  $P_r$  value, SOC, and Eqn. (6); derived the DOC concentration for the evaluation catchment (whose outlet is an observational station) by taking the area-weighted average of the local  $P_r$  DOC values from the few NHDPlus local catchments located within this catchment; 2) Derive the DOC concentration value for the evaluation catchment (whose outlet is an observational gauge) by using the average  $P_r$  value and Eqn. (6); 3) Compare the "derived" DOC concentration with the observed value at the same evaluation catchment. Note that two evaluation catchments are dropped during Step (1) for containing some NHDPlus local catchments without an effective model simulated  $P_r$ .

Figure 67 shows that our derived DOC concentration values effectively reproduce the spatial variability in the observed values. Note the unit of DOC concentration in water is mostly reported in mg/L (Schelker et al., 2012; Tian et al., 2015b; Langeveld et al., 2020). The MASE, NRMSE and  $R^2$  values are 0.73, 1.81, and 0.47, respectively, further suggesting a satisfactory performance. The scattering only occurs to a small portion of the dots, as indicated by the reddish colours. This scattering may stem from several causes, such as the limited availability of DOC observation data and the uncertainties in model development (see Section Sect. 4 for more details). Despite the scattering, the overall alignment between observed and predicted values suggests that our methods, including the generic formula and ML modelling, are appropriate and effective. The DOC evaluation performance of the SoilGrids-based model (Supplementary Fig. S13) reveals a larger systematic bias. This issue is also primarily attributed to differences in data distribution, as the  $P_r$  values in evaluation exhibit a wider range than those in training, particularly at low values (see Sect. 2.2.2). Consequently, the model struggles to predict extreme values accurately. For example, for very small  $P_r$  values in the evaluation catchments, the model tends to slightly overpredict due to the absence of such small values in the training dataset. Additionally, the typically higher SOC values in these regions further amplify the discrepancies.



**Figure 67.** Evaluation of derived DOC concentration at the catchment scale (n=3208). The solid black line indicates a 1:1 ratio. The varying colours indicate the density of points in the scatter plot.

#### 4 Uncertainty analyses

The final product, our  $P_r$  map, is subject to uncertainties from various sources. In this study, we have implemented several measures to constrain the uncertainties embedded in the input data and ML modeling exercise. We also look into the ML model parameter uncertainty via sensitivity analyses.

##### 4.1 Efforts to constrain uncertainty

###### 4.1.1 Machine learning model input data

The estimation of the DOC long-term average transformation rate,  $P_r$ , relies on SOC data from the HWSD v1.2 and SoilGrids 2.0 dataset and DOC data from the WQP stations. Despite implementing stringent catchment selection (see Section 2.2.1), the challenge of balancing data quantity and quality persists due to limited DOC measurements. Larger uncertainties in  $P_r$  are anticipated in catchments with fewer samples or those where most samples are collected in a single season. Additionally, potential uncertainties in the  $P_r$  estimation may arise from the mismatch in sampling periods between SOC and DOC datasets. It is crucial to recognize and account for these uncertainties when interpreting and using the  $P_r$  map.

525 The flowline and catchment attributes from NHDPlus constitute the primary inputs in both training and prediction phases for  
 the ML model, and thus may contribute to the uncertainty in the results. NHDPlus catchment attributes are drawn from diverse  
 sources, including remote sensing data and model simulations. Upstream-accumulated values are derived based on flowline  
 data (Wieczorek et al., 2018). A majority of attributes have been compared to equivalent variables, when available, in the  
 Geospatial Attributes of Gages for Evaluating Streamflow version II (GAGESII) dataset (Falcone et al., 2010). These  
 comparisons have demonstrated reasonably strong alignment. Inherent uncertainties may still arise from inaccurate flowline  
 530 and catchment delineation, inaccuracies in the source data, the conversion of data formats (e.g., from grid-based to catchment-  
 based), and so on. Furthermore, instances of missing data or attributes with zero-inflated values (e.g., regions highlighted in  
 white in [FigureSupplementary Fig. S5b](#)) from the NHDPlus dataset can complicate accurate data interpolation by the ML  
 model. Despite the use of the sparsity-aware technique within the XGBoost algorithm, adept at handling missing or zero-  
 inflated data to a certain extent (Chen and Guestrin, 2016), the presence of such challenges persists. Overcoming these  
 535 limitations is beyond this study's scope.

#### 4.1.2 ~~ML~~-Machine learning model development

In contrast to physical-based models with clearly pre-defined structures, ~~machine learning (ML)~~ML models endeavor to  
 discern the optimal structure from input data through the training process. Consequently, uncertainty may emerge at any stage  
 of model development, as detailed in [SectionSect. 2.3](#). To mitigate model uncertainty, we employ well-established strategies  
 540 prevalent in diverse applications (Abeshu et al., 2022; Delavar et al., 2019; Li et al., 2022). These encompass techniques such  
 as transformation of input data, training and testing splits, feature selection, hyperparameter tuning, and cross-validation (refer  
 to previous sections for details). These measures aim to constrain the uncertainties inherent in model development processes  
 and fortify the model's predictive capabilities, for example by refining the interpretability of input data, mitigating the risk of  
 overfitting, enhancing generalization performance, and minimizing the introduction of potentially noisy predictors.

545 In addition to the commonly adopted strategies in using XGBoost and the other ML techniques, we augment the control of  
 model uncertainty through a representativeness check. This check ensures alignment between the distribution of model  
 parameters used during training and those applied in predictions. This additional step serves to enhance the model's  
 transferability from the training catchment to the broader CONUS domain. To gauge the representativeness of our chosen  
 550 predictors, we conducted a Cumulative Distribution Function (CDF) comparison for each parameter between the observational  
 dataset (derived from 2595 independent catchments) and the entire CONUS dataset (comprising approximately 2.6 million  
 local catchments in NHDPlus). For this comparison, we assess the relative difference in the 5th, 25th, 50th, 75th, and 95th  
 percentiles between the two CDFs. As an illustration, the relative difference for the 5th percentile is computed as the ratio of  
 the difference between the 5th percentile of the available  $P_r$  data and that of the entire CONUS data to their average. Table [42](#)  
 555 provides a summary of the CDF comparison of the 15 selected predictors ([also-see-supplementary-FigureSupplementary Fig.](#)

S6). A predictor is deemed representative of the whole CONUS if the average relative difference is less than 0.75. Following Abeshu et al. (2022), the choice of the 0.75 threshold strikes a balance between maintaining data representativeness and avoiding the exclusion of too many predictors. Three predictors, namely "basin-area", "per\_hwetlandBASIN\_AREA", "NLCD01\_95", and "per\_shrubNLCD01\_52", have failed the representativeness check and are consequently excluded. Note that the ML model performance has only slightly changed after reducing the number of predictors from 15 to 12, as shown in the supplementary Figure S7. Following the same process, the SoilGrids-based model excludes "NLCD01\_95" during the representativeness check, resulting in 12 out of 13 predictors being retained for the final optimal model (Supplementary Table S5).

Table 42. Representativeness of XGBoost model input predictors over CONUS.

Attributes	Relative difference in percentiles between $P_i$ -available and whole_conus data					Average
	5th	25th	50th	75th	95th	
basin-areaBASIN_A REA	1.941	1.728	1.669	1.794	1.900	1.806
per_hwetlandNLCD0 1_95	0.667	0.667	0.842	1.144	1.529	0.969
per_shrubNLCD01_5 2	0.353	0.624	1.224	1.482	0.889	0.914
per_canopyCNPY11 BUFF100	1.684	1.090	0.427	0.080	0.078	0.672
per_wetlandNLCD 01_90	0.769	0.314	0.461	0.621	0.807	0.594
per_forestNLCD01 42	0.667	0.559	0.651	0.502	0.225	0.521
elev_related	0.769	0.806	0.320	0.621	0.008	0.505
hydro_related	0.584	0.898	0.316	0.108	0.106	0.402
per_soilmsHGBD	0.955	0.264	0.152	0.095	0.255	0.344
index_tqsubCONTA CT	0.166	0.135	0.248	0.292	0.393	0.247
index_bflowBFI	0.476	0.304	0.152	0.002	0.027	0.192
per_humidityRH	0.197	0.103	0.015	0.014	0.014	0.068
soil_texture_related	0.095	0.071	0.068	0.071	0.015	0.064
ave_wetdayCWD	0.063	0.065	0.028	0.053	0.033	0.048
temp_related	0.035	0.034	0.009	0.029	0.006	0.023

Abbreviations: BASIN\_AREA (catchment area); NLCD01\_95 (areal percentage of herbaceous wetlands); NLCD01\_52 (areal percentage of shrub); CNPY11 BUFF100 (areal percentage of canopy in the riparian buffer); NLCD01\_90 (areal percentage of woody wetlands); NLCD01\_42 (areal percentage of evergreen forest); elev\_related (merged predictor for mean/min/max elevation); hydro\_related (merged predictor representing recharge, runoff, and precipitation); HGBD (areal percentage of Hydrologic Group BD soil); CONTACT (subsurface contact time); BFI (base flow index); RH (relative humidity); soil\_texture\_related (merged predictor for silt and sand content); CWD (consecutive wet days); temp\_related (merged predictor encompassing potential evapotranspiration, first/last freeze timing, snow fraction, actual evapotranspiration, and mean/min/max temperature); For detailed descriptions, refer to Supplementary Tables S2 and S3.

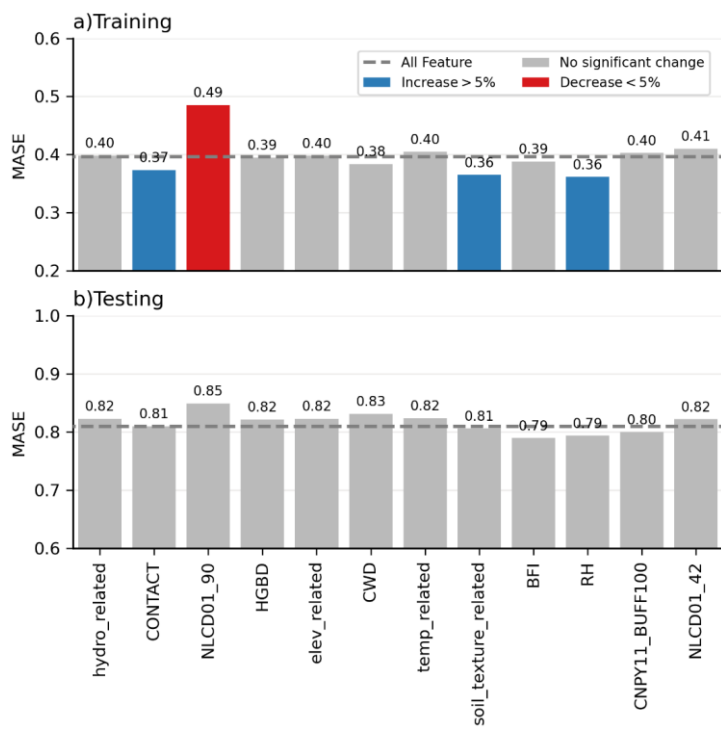
Formatted Table

Formatted: English (United Kingdom)

4.2 Sensitivity analyses

Model sensitivity analysis (SA) involves probing the importance of uncertainties in model parameters (Loucks and Van Beek, 2017). We examine our model's sensitivity to each selected predictor using two different methods: 1) dropping one predictor

at a time and tracking the changes in model performance, and 2) the Sobol sensitivity analysis approach (Sobol, 2001). Figure 78 demonstrates the model performance difference in the training and testing phases after dropping one of the 12 variables. Blue, red, and grey colors are employed to indicate whether dropping the corresponding predictor will result in an increase, decrease, or insignificant change in the model's performance, respectively. A 5% threshold is chosen to determine the significance of the change. In general, the shifting pattern in MASE scores remains consistent between the training and testing phases. However, the alterations in MASE values for most predictors, particularly during the testing phase, are minimal or even negligible. In other words, the model appears to be insensitive to most predictors according to this first sensitivity analysis method.



**Figure 78.** Sensitivity of XGBoost model to predictors in the training and testing phases. The MASE value is represented by the blue, red, and grey bars, indicating whether the model performance increases, decreases, or remains relatively unchanged after dropping the corresponding predictor. The dashed grey line indicates the model performance with all variables included. Abbreviations: hydro\_related (merged predictor representing recharge, runoff, and precipitation); CONTACT (subsurface contact time); NLCD01\_90 (areal percentage of woody wetlands); HGBD (areal percentage of Hydrologic Group BD soil); elev\_related (merged

595 predictor for mean/min/max elevation); CWD (consecutive wet days); temp\_related (merged predictor encompassing potential  
evapotranspiration, first/last freeze timing, snow fraction, actual evapotranspiration, and mean/min/max temperature);  
600 soil\_texture\_related (merged predictor for silt and sand content); BFI (base flow index); RH (relative humidity); CNPY11\_BUFF100  
(areal percentage of canopy in the riparian buffer); NLCD01\_42 (areal percentage of evergreen forest). For detailed descriptions,  
refer to Supplementary Tables S2 and S3.

The Sobol sensitivity analysis is a widely used variance-based global sensitivity analysis method (Borgonovo and Plischke, 2016). It provides two indices: First-order Index (S1), which measures the sensitivity of an individual predictor itself (local  
600 variance), and Total Index (ST), which accounts for the effects of both an individual predictor itself and its interactions with  
any other predictors (global variance) (Saltelli, 2002; Saltelli et al., 2010). These interactions, which can be of any order, can  
be isolated. For instance, second and higher-order interactions can be isolated by subtracting S1 from ST. The results from the  
Sobol test are summarized in Table 53. The distribution of S1 is highly right-skewed, suggesting that the model exhibits  
insensitivity to most predictors if only local variance is considered. There are, however, a few exceptions, such as  
605 "hydro\_related", and "temp\_related", which present high S1 values. The global variance, represented by the ST index, paints  
a somewhat different picture. When considering the ST index, a broad set of predictors emerge as sensitive, particularly those  
with ST values exceeding 0.1. It's worth noting that these predictors also hold high rankings in the predictor selection, as shown  
in Table 2 Fig. 3. Furthermore, it is significant that 11 out of the total 12 predictors show a normalized difference between S1  
and ST (calculated as (ST-S1)/ST) greater than 50%. This observation underscores the significant interactions among the  
610 predictors (Saltelli et al., 2010). This suggests that if a predictor is dropped, the remaining predictors could potentially  
compensate for its absence, highlighting the nonlinear, high-order interdependence among the predictors in our model.

**Table 53. Sobol sensitivity analysis results for the 12 selected predictors.**

Predictors	Total Indices (ST)	First Order Indices (S1)	Difference ((ST-S1)/ST)
hydro_related	0.466	0.291	0.375
temp_related	0.311	0.141	0.546
ave_wetdayCWD	0.207	0.044	0.788
index_tqsubCONT	0.143	0.003	0.977
ACT			
per_canopyCNPY1	0.132	0.028	0.787
1_BUFF100			
per_wetlandNLC	0.125	0.049	0.608
D01_90			
elev_related	0.087	0.017	0.806
index_bflowBFI	0.072	0.012	0.831
per_rhumidityRH	0.062	0.010	0.836
soil_texture_related	0.034	0.000	1.000
per_eforestNLCD01	0.024	0.005	0.798
42			
per_soilmsHGBD	0.013	0.002	0.873

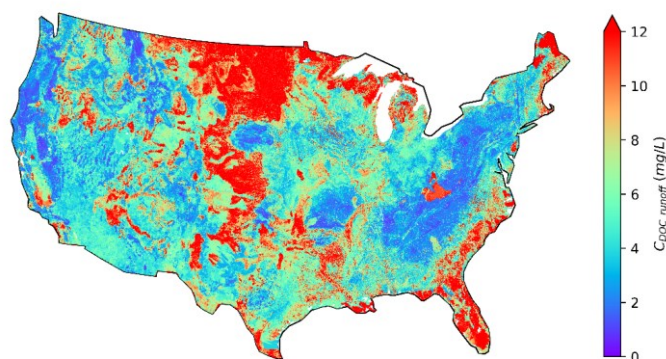


Abbreviations: hydro\_related (merged predictor representing recharge, runoff, and precipitation); temp\_related (merged predictor encompassing potential evapotranspiration, first/last freeze timing, snow fraction, actual evapotranspiration, and mean/min/max temperature); CWD (consecutive wet days); CONTACT (subsurface contact time); CNPY11\_BUFF100 (areal percentage of canopy in the riparian buffer); NLCD01\_90 (areal percentage of woody wetlands); elev\_related (merged predictor for mean/min/max elevation); BFI (base flow index); RH (relative humidity); soil\_texture\_related (merged predictor for silt and sand content); NLCD01\_42 (areal percentage of evergreen forest); HGBD (areal percentage of Hydrologic Group BD soil); For detailed descriptions, refer to Supplementary Tables S2 and S3.

The above sensitivity analyses suggest that our model exhibits low sensitivity to most predictors when considering their individual (local) impact. However, the Sobol sensitivity analysis uncovers a heightened degree of sensitivity in the context of global effects, particularly given the significant interactions among the predictors. A similar sensitivity analysis was conducted for the SoilGrids-based model, yielding the same conclusions (Supplementary Fig. S14 and Supplementary Table S6).

## 5 Potential use and limitations

The  $P_r$  map has several promising uses. For instance, one of the pivotal applications of the  $P_r$  map is to estimate the lateral leaching of DOC. Figure 89, as an illustration, shows a  $C_{DOC\_runoff}$  map over CONUS depicting the long-term average concentration of DOC in the leaching flux at over two million NHDPlus local catchments. This map is derived based on Eqn. (4), leveraging the  $P_r$  map in Fig. 6 and the top-layer SOC data from HWDS1.2. Due to missing data in the HWSD 1km SOC map at about 0.6 million NHDPlus local catchments, we cannot calculate the  $C_{DOC\_runoff}$  values over those catchments.



**Figure 89.** Calculated CONUS map of DOC concentration in leaching flux from soils to headwater streams over 2.6 million NHDPlus flowlines.

The spatial patterns of the  $C_{DOC\_runoff}$  map are highly correlated to those of the  $P_r$  (see Figure 5 Fig 6) and SOC map (see supplementary Figure Supplementary Fig S5a). Notably, the  $C_{DOC\_runoff}$  values are high in regions with extremely high SOC values. Additionally, the  $C_{DOC\_runoff}$  values are high in North Dakota, Montana, and southern coasts, where the  $P_r$  values are high. Interestingly, the influences of  $P_r$  and SOC can counterbalance each other in some places. For instance, in the upper Rocky Mountains, the SOC storage is abundant due to the presence of forests. However, the low temperature in this region hinders microbial activities, resulting in extremely low  $P_r$  values. As a result, the concentration of DOC leaching flux is relatively low. Moreover, the spatial coverage of wetlands also appears to be relevant (see supplementary Figure Supplementary Fig. S5b), which is consistent with the suggested crucial role of wetlands in riverine DOC dynamics (Duan et al., 2017; Leibowitz et al., 2023). For instance, high  $C_{DOC\_runoff}$  values are observed in upper Minnesota, Florida, and Louisiana, where wetlands are prevalent. In places with few wetlands, like Nevada, Arizona, and New Mexico, the leaching flux concentration is considerably lower.

There are at least two other potential uses of the  $P_r$  map: 1) It can support large-scale DOC modeling over CONUS or a major river basin. For instance, testing the use of the map within the framework of the Energy Exascale Earth System Model (Burrows-Golaz et al., 2019; Caldwell et al., 2019; Golaz et al., 2019; Burrows et al., 2020) is ongoing and will be reported in the near future. 2) It can be used to provide a quick estimation of riverine DOC concentration or flux at any headwater catchments where no DOC observations are available.

We caution the potential users of the  $P_r$  map with several limitations in the methods invoked. Firstly, the  $P_r$  values in the map account for the spatial heterogeneity of various DOC-related processes and factors only in a long-term average sense owing to the limited data availability, i.e., the SOC reanalysis data are long-term averages, and the observed riverine DOC data are only available at irregular time intervals. While we believe that such a  $P_r$  map is a critical step in effectively capturing the spatial heterogeneity of the relevant processes and environmental factors, incorporating their temporal dynamics is beyond the scope of this study and left for future work. Secondly, the ML techniques are not process-based and thus do not yet offer rich insight into the relevant mechanisms. To improve our understanding of the DOC-related processes, the  $P_r$  map should be used in conjunction with other observational data, process-based models, and carefully designed numerical experiments. Third, the lack of direct measurements of  $P_r$  necessitates the use of indirect validation methods. To further enhance robustness, we encourage the design and implementation of new field experiments guided by our lumped parameter approach. Last but not least, the ML model has been trained with the data in the CONUS domain only, so it may not be transferable beyond CONUS.

Our lumped parameter approach and machine learning-based parameterization strategy are designed to generalize beyond the CONUS and scale globally. The framework is inherently generic, independent of site-specific characteristics, and supported by machine learning techniques adaptable to diverse regions. The CONUS study area, characterized by substantial spatial heterogeneity, provides a robust foundation for demonstrating this generalizability. However, extending the framework to a

global scale introduces challenges, particularly in data availability and variability in environmental conditions. Addressing these requires extensive observational data collection, especially riverine DOC observations, leveraging public datasets, literature, and increased fieldwork for enhanced coverage. At the global scale, managing increased uncertainties is crucial, as larger variability is expected compared to the CONUS-based parameterization. Efforts should focus on assembling comprehensive catchment attributes while maintaining flexibility in their significance assessment, allowing the machine learning model to determine their importance contextually. High-priority attributes identified in this study (Fig. 3), such as woody wetland percentage, should receive particular attention as they are likely critical in other regions.

## 6 Data and code availability

The resulting  $P_r$  and  $C_{DOC\_runoff}$  maps over CONUS are freely available at <https://doi.org/10.5281/zenodo.8339372> (Li et al., 2024). <https://zenodo.org/records/14563816> (Li et al., 2024). The Zenodo repository includes the following resources: a) Pr.gpkg – a 9.9 GB GeoPackage file containing data on  $P_r$ , SOC, and DOC, derived using SOC data from HWSD v1.2 and SoilGrids 2.0 across over 2.6 million NHDPlus local catchments. This file also includes COMID and local catchment boundary polygons and is compatible with GIS software such as QGIS, ArcGIS, and Python libraries like GeoPandas for analysis and editing; b) PNG images – two high-resolution PNG files illustrating the HWSD-based and SoilGrids-based model-simulated  $P_r$  maps across over 2.6 million NHDPlus local catchments; c) Required input files – files necessary to reproduce the reported results; and d) ReadMe document – a text file providing detailed descriptions of each resource in the Zenodo repository. The input data are obtained from the water quality portal (<https://www.waterqualitydata.us/>), NHDPlus (<https://www.epa.gov/waterdata/nhdplus-national-data>), ScienceBase (<https://doi.org/10.5066/F7765D7V>) and HWSD v1.2 (<https://www.fao.org/soils-portal/data-hub/soil-maps-and-databases/harmonized-world-soil-database-v12/en/>) and SoilGrids2.0 (<https://files.isric.org/soilgrids/latest/data/>). Additionally, the Python scripts used for feature selection, model training, and evaluation are available on the Github repository at <https://github.com/Ceyxleo/DOC-Param-Map>.

## 7 Conclusions

We developed two new maps of  $P_r$ , the transformation rate from SOC concentration in soil to DOC concentration in the leaching flux, over CONUS, based on SOC data from the HWSD v1.2 and SoilGrids 2.0. Evaluation of derived DOC concentrations at over 3000 headwater WQP stations confirms the robustness of our methodology, including which incorporates a generic formula linking SOC and DOC via  $P_r$ , riverine DOC observations, environmental variables, and the ML techniques that effectively capture high-order, nonlinear relationships between  $P_r$  and the environmental variables. Such a map did not exist before and is. These  $P_r$  maps, the first of their kind, are highly valuable for large-scale DOC modeling and for improving our understanding of the DOC-related processes across the land-river continuum.

700 LL performed the analysis with the inputs from the co-authors, prepared the figures, and wrote the first draft. HL devised the conceptual idea and supervised the study. GA provided frequent assistance in processing the data and developing the model. All the co-authors contributed to the writing.

**Competing interests**

At least one of the (co-)authors is a member of the editorial board of Earth System Science Data.

705 **Disclaimer**

Publisher's note: Copernicus Publications remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Acknowledgments**

710 This research is supported by the Office of Science of the U.S. Department of Energy [Biological and Environmental Research](#) as part of the Earth System Model Development program area through the Energy Exascale Earth System Model (E3SM) project. The Pacific Northwest National Laboratory is operated by Battelle for the U.S. Department of Energy under Contract DE-AC05-76RL01830.

- Abeshu, G. W., Li, H.-Y., Zhu, Z., Tan, Z., and Leung, L. R.: Median bed-material sediment particle size across rivers in the contiguous US, *Earth Syst Sci Data*, 14, 929–942, <https://doi.org/10.5194/essd-14-929-2022>, 2022.
- Afan, H. A., El-shafie, A., Mohtar, W. H. M. W., and Yaseen, Z. M.: Past, present and prospect of an Artificial Intelligence (AI) based model for sediment transport prediction, *J Hydrol (Amst)*, 541, 902–913, <https://doi.org/10.1016/j.jhydrol.2016.07.048>, 2016.
- Akiba, T., Sano, S., Yanase, T., Ohta, T., and Koyama, M.: Optuna: A Next-generation Hyperparameter Optimization Framework, in: *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2623–2631, <https://doi.org/10.1145/3292500.3330701>, 2019.
- Alebachew, M. A., Ye, S., Li, H., Huang, M., Leung, L. R., Fiori, A., and Sivapalan, M.: Regionalization of subsurface stormflow parameters of hydrologic models: Up-scaling from physically based numerical simulations at hillslope scale, *J Hydrol (Amst)*, 519, 683–698, <https://doi.org/10.1016/j.jhydrol.2014.07.018>, 2014.
- Autio, I., Soinne, H., Helin, J., Asmala, E., and Hoikkala, L.: Effect of catchment land use and soil type on the concentration, quality, and bacterial degradation of riverine dissolved organic matter, *Ambio*, 45, 331–349, <https://doi.org/10.1007/s13280-015-0724-y>, 2016.
- Ayata, S.-D., Irissou, J.-O., Aubert, A., Berline, L., Dutay, J.-C., Mayot, N., Nieblas, A.-E., ~~D’Ortenzio~~D’Ortenzio, F., Palmiéri, J., Reygondeau, G., Rossi, V., and Guieu, C.: Regionalisation of the Mediterranean basin, a MERMEX synthesis, *Prog Oceanogr*, 163, 7–20, <https://doi.org/10.1016/j.pocean.2017.09.016>, 2018.
- Batjes, N. H.: A taxotransfer rule-based approach for filling gaps in measured soil data in primary SOTER databases (Version 1.1) Global Environment Facility United Nations Environment Programme Netherlands Ministry of Housing, Spatial Planning and the Environment, 2003.
- Baum, A., Rixen, T., and Samiaji, J.: Relevance of peat draining rivers in central Sumatra for the riverine input of dissolved organic carbon into the ocean, *Estuar Coast Shelf Sci*, 73, 563–570, <https://doi.org/10.1016/j.ecss.2007.02.012>, 2007.
- Best, M. J., Pryor, M., Clark, D. B., Rooney, G. G., Essery, R. . L. H., Ménard, C. B., Edwards, J. M., Hendry, M. A., Porson, A., Gedney, N., Mercado, L. M., Sitch, S., Blyth, E., Boucher, O., Cox, P. M., Grimmond, C. S. B., and Harding, R. J.: The Joint UK Land Environment Simulator (JULES), model description – Part 1: Energy and water fluxes, *Geosci Model Dev*, 4, 677–699, <https://doi.org/10.5194/gmd-4-677-2011>, 2011.
- Borgonovo, E. and Plischke, E.: Sensitivity analysis: A review of recent advances, *Eur J Oper Res*, 248, 869–887, <https://doi.org/10.1016/j.ejor.2015.06.032>, 2016.
- Brooks, P. D., McKnight, D. M., and Bencala, K. E.: The relationship between soil heterotrophic activity, soil dissolved organic carbon (DOC) leachate, and catchment-scale DOC export in headwater catchments, *Water Resour Res*, 35, 1895–1902, <https://doi.org/10.1029/1998WR900125>, 1999.

- 745 Burrows, S. M., Maltrud, M., Yang, X., Zhu, Q., Jeffery, N., Shi, X., Ricciuto, D., Wang, S., Bisht, G., Tang, J., Wolfe, J., Harrop, B. E., Singh, B., Brent, L., Baldwin, S., Zhou, T., Cameron-Smith, P., Keen, N., Collier, N., Xu, M., Hunke, E. C., Elliott, S. M., Turner, A. K., Li, H., Wang, H., Golaz, J. -C., Bond-Lamberty, B., Hoffinan, F. M., Riley, W. J., Thornton, P. E., Calvin, K., and Leung, L. R.: The DOE E3SM v1.1 Biogeochemistry Configuration: Description and Simulated Ecosystem-Climate Responses to Historical Changes in Forcing, *J Adv Model Earth Syst*, 12, <https://doi.org/10.1029/2019MS001766>,  
750 2020.
- Caldwell, P. M., Mametjanov, A., Tang, Q., Van Roekel, L. P., Golaz, J. C., Lin, W., Bader, D. C., Keen, N. D., Feng, Y., Jacob, R., Maltrud, M. E., Roberts, A. F., Taylor, M. A., Veneziani, M., Wang, H., Wolfe, J. D., Balaguru, K., Cameron-Smith, P., Dong, L., Klein, S. A., Leung, L. R., Li, H. Y., Li, Q., Liu, X., Neale, R. B., Pinheiro, M., Qian, Y., Ullrich, P. A., Xie, S., Yang, Y., Zhang, Y., Zhang, K., and Zhou, T.: The DOE E3SM Coupled Model Version 1: Description and Results at High  
755 Resolution, *J Adv Model Earth Syst*, 11, 4095–4146, <https://doi.org/10.1029/2019MS001870>, 2019.
- Camino-Serrano, M., Gielen, B., Luyssaert, S., Ciais, P., Vicca, S., Guenet, B., Vos, B. De, Cools, N., Ahrens, B., Altaf Arain, M., Borken, W., Clarke, N., Clarkson, B., Cummins, T., Don, A., Pannatier, E. G., Laudon, H., Moore, T., Nieminen, T. M., Nilsson, M. B., Peichl, M., Schwendenmann, L., Siemens, J., and Janssens, I.: Linking variability in soil solution dissolved organic carbon to climate, soil type, and vegetation type, *Global Biogeochem Cycles*, 28, 497–509,  
760 <https://doi.org/10.1002/2013GB004726>, 2014.
- Chegini, T., Li, H.-Y., and Leung, L.: HyRiver: Hydroclimate Data Retriever, *J Open Source Softw*, 6, 3175, <https://doi.org/10.21105/joss.03175>, 2021.
- Chen, T. and Guestrin, C.: XGBoost: A scalable tree boosting system, in: *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 785–794, <https://doi.org/10.1145/2939672.2939785>, 2016.
- 765 Chow, V. Te, Maidment, D. R., and Mays, L. W.: *Applied hydrology*, McGraw-Hill, 572 pp., 1988.
- [Dai, Y., Shangguan, W., Wei, N., Xin, Q., Yuan, H., Zhang, S., Liu, S., Lu, X., Wang, D., and Yan, F.: A review of the global soil property maps for Earth system models, \*SOIL\*, 5, 137–158, <https://doi.org/10.5194/soil-5-137-2019>, 2019.](#)
- Daoud, J. I.: Multicollinearity and Regression Analysis, in: *Journal of Physics: Conference Series*, <https://doi.org/10.1088/1742-6596/949/1/012009>, 2018.
- 770 [Davidson, E. A. and Janssens, I. A.: Temperature sensitivity of soil carbon decomposition and feedbacks to climate change, \*https://doi.org/10.1038/nature04514\*, 9 March 2006.](#)
- Delavar, M. R., Gholami, A., Shiran, G. R., Rashidi, Y., Nakhaeizadeh, G. R., Fedra, K., and Afshar, S. H.: A novel method for improving air pollution prediction based on machine learning approaches: A case study applied to the capital city of Tehran, *ISPRS Int J Geoinf*, 8, <https://doi.org/10.3390/ijgi8020099>, 2019.
- 775 Doron, M., Brasseur, P., and Brankart, J. M.: Stochastic estimation of biogeochemical parameters of a 3D ocean coupled physical-biogeochemical model: Twin experiments, *Journal of Marine Systems*, 87, 194–207, <https://doi.org/10.1016/j.jmarsys.2011.04.001>, 2011.

Duan, S., He, Y., Kaushal, S. S., Bianchi, T. S., Ward, N. D., and Guo, L.: Impact of Wetland Decline on Decreasing Dissolved Organic Carbon Concentrations along the Mississippi River Continuum, *Front Mar Sci*, 3, <https://doi.org/10.3389/fmars.2016.00280>, 2017.

Duarte, C. M.: Reviews and syntheses: Hidden forests, the role of vegetated coastal habitats in the ocean carbon budget, <https://doi.org/10.5194/bg-14-301-2017>, 23 January 2017.

Ducharme, A., Golaz, C., Leblois, E., Laval, K., Polcher, J., Ledoux, E., and De Marsily, G.: Development of a high resolution runoff routing model, calibration and application to assess runoff from the LMD GCM, *J Hydrol (Amst)*, 280, 207–228, [https://doi.org/10.1016/S0022-1694\(03\)00230-0](https://doi.org/10.1016/S0022-1694(03)00230-0), 2003.

Dupas, R., Curie, F., Gascuel-Oudou, C., Moatar, F., Delmas, M., Parnaudeau, V., and Durand, P.: Assessing N emissions in surface water at the national level: Comparison of country-wide vs. regionalized models, *Science of the Total Environment*, 443, 152–162, <https://doi.org/10.1016/j.scitotenv.2012.10.011>, 2013.

Falcone, J. A., Carlisle, D. M., Wolock, D. M., and Meador, M. R.: GAGES: A stream gage database for evaluating natural and altered flow conditions in the conterminous United States, *Data Papers Ecology*, 621 pp., 2010.

Fan, C., Song, C., Liu, K., Ke, L., Xue, B., Chen, T., Fu, C., and Cheng, J.: Century-Scale Reconstruction of Water Storage Changes of the Largest Lake in the Inner Mongolia Plateau Using a Machine Learning Approach, *Water Resour Res*, 57, <https://doi.org/10.1029/2020WR028831>, 2021.

Finlay, J., Neff, J., Zimov, S., Davydova, A., and Davydov, S.: Snowmelt dominance of dissolved organic carbon in high-latitude watersheds: Implications for characterization and flux of river DOC, *Geophys Res Lett*, 33, <https://doi.org/10.1029/2006GL025754>, 2006.

Fischer, G., Nachtergaele, F., Prieler, S., Van Velthuisen, H. T., Verelst, L., and Wiberg, D.: Global Agro-ecological Zones Assessment for Agriculture (GAEZ 2008), Laxenburg, Austria and FAO, Rome, Italy, 2008.

Futter, M. N., Butterfield, D., Cosby, B. J., Dillon, P. J., Wade, A. J., and Whitehead, P. G.: Modeling the mechanisms that control in-stream dissolved organic carbon dynamics in upland and forested catchments, *Water Resour Res*, 43, <https://doi.org/10.1029/2006WR004960>, 2007.

Golaz, J. C., Caldwell, P. M., Van Roekel, L. P., Petersen, M. R., Tang, Q., Wolfe, J. D., Abeshu, G., Anantharaj, V., Asay-Davis, X. S., Bader, D. C., Baldwin, S. A., Bisht, G., Bogenschutz, P. A., Branstetter, M., Brunke, M. A., Brus, S. R., Burrows, S. M., Cameron-Smith, P. J., Donahue, A. S., Deakin, M., Easter, R. C., Evans, K. J., Feng, Y., Flanner, M., Foucar, J. G., Fyke, J. G., Griffin, B. M., Hannay, C., Harrop, B. E., Hoffman, M. J., Hunke, E. C., Jacob, R. L., Jacobsen, D. W., Jeffery, N., Jones, P. W., Keen, N. D., Klein, S. A., Larson, V. E., Leung, L. R., Li, H. Y., Lin, W., Lipscomb, W. H., Ma, P. L., Mahajan, S., Maltrud, M. E., Mamatjanov, A., McClean, J. L., McCoy, R. B., Neale, R. B., Price, S. F., Qian, Y., Rasch, P. J., Reeves Eyre, J. E. J., Riley, W. J., Ringler, T. D., Roberts, A. F., Roesler, E. L., Salinger, A. G., Shaheen, Z., Shi, X., Singh, B., Tang, J., Taylor, M. A., Thornton, P. E., Turner, A. K., Veneziani, M., Wan, H., Wang, H., Wang, S., Williams, D. N., Wolfram, P. J., Worley, P. H., Xie, S., Yang, Y., Yoon, J. H., Zelinka, M. D., Zender, C. S., Zeng, X., Zhang, C., Zhang, K.,

- Zhang, Y., Zheng, X., Zhou, T., and Zhu, Q.: The DOE E3SM Coupled Model Version 1: Overview and Evaluation at Standard Resolution, *J Adv Model Earth Syst*, 11, 2089–2129, <https://doi.org/10.1029/2018MS001603>, 2019.
- Gupta, H. V., Kling, H., Yilmaz, K. K., and Martinez, G. F.: Decomposition of the mean squared error and NSE performance criteria: Implications for improving hydrological modelling, *J Hydrol (Amst)*, 377, 80–91, <https://doi.org/10.1016/j.jhydrol.2009.08.003>, 2009.
- Han, X., Franssen, H. J., Hamby, D. M.: A Comparison of Sensitivity Analysis Techniques, *Health Phys*, 68, 195–204, <https://doi.org/10.1097/00004032-199502000-00005>, 1995.
- Hansell, D., Carlson, C., Repeta, D., and Schlitzer, R.: Dissolved Organic Matter in the Ocean: A Controversy Stimulates New Insights, *Oceanography*, 22, 202–211, <https://doi.org/10.5670/oceanog.2009.109>, 2009.
- He, C., Yang, C. J., Turowski, J. M., Ott, R. F., Braun, J., Tang, H., Ghantous, S., Yuan, X., and De Quay, G. S.: A global dataset of the shape of drainage systems, *Earth Syst Sci Data*, 16, 1151–1166, <https://doi.org/10.5194/essd-16-1151-2024>, 2024.
- Helton, A. M., Wright, M. S., Bernhardt, E. S., Poole, G. C., Cory, R. M., and Stanford, J. A.: Dissolved organic carbon lability increases with water residence time in the alluvial aquifer of a river floodplain ecosystem, *J Geophys Res Biogeosci*, 120, 693–706, <https://doi.org/10.1002/2014JG002832>, 2015.
- Hengl, T., De Jesus, J. M., Heuvelink, G. B. M., Gonzalez, M. R., Kilibarda, M., Blagotić, A., Shangguan, W., Wright, M. N., Geng, X., Bauer-Marschallinger, B., Guevara, M. A., Vargas, R., MacMillan, R. A., Batjes, N. H., Leenaars, J. G. B., Ribeiro, E., Wheeler, I., Mantel, S., and Kempen, B.: SoilGrids250m: Global gridded soil information based on machine learning, *PLoS One*, 12, <https://doi.org/10.1371/journal.pone.0169748>, 2017.
- Hyndman, R. J. and Koehler, A. B.: Another look at measures of forecast accuracy, *Int J Forecast*, 22, 679–688, <https://doi.org/10.1016/j.ijforecast.2006.03.001>, 2006.
- Jing, X., Tian, G., Li, M., and Javeed, S. A.: Research on the spatial and temporal differences of ~~ehina's~~china's provincial carbon emissions and ecological compensation based on land carbon budget accounting, *Int J Environ Res Public Health*, 18, <https://doi.org/10.3390/ijerph182412892>, 2021.
- Kaiser, K. and Kalbitz, K.: Cycling downwards - dissolved organic matter in soils, *Soil Biol Biochem*, 52, 29–32, <https://doi.org/10.1016/j.soilbio.2012.04.002>, 2012.
- Kalbitz, K., Solinger, S., Park, J.-H., Michalzik, B., and Matzner, E.: CONTROLS ON THE DYNAMICS OF DISSOLVED ORGANIC MATTER IN SOILS: A REVIEW, *Soil Sci*, 165, 277–304, <https://doi.org/10.1097/00010694-200004000-00001>, 2000.
- Knoben, W. J. M., Freer, J. E., and Woods, R. A.: Technical note: Inherent benchmark or not? Comparing Nash-Sutcliffe and Kling-Gupta efficiency scores, *Hydrol Earth Syst Sci*, 23, 4323–4331, <https://doi.org/10.5194/hess-23-4323-2019>, 2019.
- Kortelainen, N. M. and Karhu, J. A.: Tracing the decomposition of dissolved organic carbon in artificial groundwater recharge using carbon isotope ratios, *Applied Geochemistry*, 21, 547–562, <https://doi.org/10.1016/j.apgeochem.2006.01.004>, 2006.



[illegible]

- Christensen, J. R., Bellmore, R. A., and Lane, C. R.: Atmospheric connectivity classification links wetlands with stream water quality, *Nature Water*, 1, 370–380, <https://doi.org/10.1038/s44221-023-00057-w>, 2023.
- Li, H. Y. and Sivapalan, M.: Functional approach to exploring climatic and landscape controls on runoff generation: 2 Timing of runoff storm response, *Water Resour Res*, 50, 9323–9342, <https://doi.org/10.1002/2014WR016308>, 2014.
- Li, H. Y., Sivapalan, M., Tian, F., and Harman, C.: Functional approach to exploring climatic and landscape controls of runoff generation: 1. Behavioral constraints on runoff volume, *Water Resour Res*, 50, 9300–9322, <https://doi.org/10.1002/2014WR016307>, 2014.
- Li, H., Wigmosta, M. S., Wu, H., Huang, M., Ke, Y., Coleman, A. M., and Leung, L. R.: A physically based runoff routing model for land surface and earth system models, *J Hydrometeorol*, 14, 808–828, <https://doi.org/10.1175/JHM-D-12-015.1>, 2013.
- Li, L., Li, H.-Y., and Abeshu, G., Tang, J., Leung, L. R., Liao, C., Tan, Z., Tian, W., H., Thornton, P., & Yang, X.: Deriving a Transformation Rate MapMaps of Dissolved Organic Carbon overin the Contiguous U.S., Zenodo [Data set], <https://doi.org/10.5281/zenodo.8339372>, <https://doi.org/10.5281/zenodo.14563816>, 2024.
- Li, L., Qiao, J., Yu, G., Wang, L., Li, H. Y., Liao, C., and Zhu, Z.: Interpretable tree-based ensemble model for predicting beach water quality, *Water Res*, 211, <https://doi.org/10.1016/j.watres.2022.118078>, 2022.
- Li, M., Peng, C., Zhou, X., Yang, Y., Guo, Y., Shi, G., and Zhu, Q.: Modeling Global Riverine DOC Flux Dynamics From 1951 to 2015, *J Adv Model Earth Syst*, 11, 514–530, <https://doi.org/10.1029/2018MS001363>, 2019.
- Liao, C., Zhuang, Q., Leung, L. R., and Guo, L.: Quantifying Dissolved Organic Carbon Dynamics Using a Three-Dimensional Terrestrial Ecosystem Model at High Spatial-Temporal Resolutions, *J Adv Model Earth Syst*, 11, 4489–4512, <https://doi.org/10.1029/2019MS001792>, 2019.
- Liu, H., Li, Q., Bai, Y., Yang, C., Wang, J., Zhou, Q., Hu, S., Shi, T., Liao, X., and Wu, G.: Improving satellite retrieval of oceanic particulate organic carbon concentrations using machine learning methods, *Remote Sens Environ*, 256, <https://doi.org/10.1016/j.rse.2021.112316>, 2021.
- Lønborg, C., Carreira, C., Jickells, T., and Álvarez-Salgado, X. A.: Impacts of Global Change on Ocean Dissolved Organic Carbon (DOC) Cycling, <https://doi.org/10.3389/fmars.2020.00466>, 23 June 2020.
- Loucks, D.P. and Van Beek, E.: *Water Resource Systems Planning and Management: An Introduction to Methods, Models, and Applications*, Springer International Publishing, 624 pp., ISBN 9783319442341, 2017.
- Ludwig, W., Probst, J.-L., and Kempe, S.: Predicting the oceanic input of organic carbon by continental erosion, *Global Biogeochem Cycles*, 10, 23–41, <https://doi.org/10.1029/95GB02925>, 1996.
- Lundberg, S. M. and Lee, S.-I.: A Unified Approach to Interpreting Model Predictions, in: *Advances in Neural Information Processing Systems*, 2017.
- McBratney, A. B., Mendonça Santos, M. L., and Minasny, B.: On digital soil mapping, *Geoderma*, 117, 3–52, [https://doi.org/10.1016/S0016-7061\(03\)00223-4](https://doi.org/10.1016/S0016-7061(03)00223-4), 2003.

- McKay, L.; Bondelid, T.; Dewald, T.; Johnston, J.; Moore, R.; Rea, A., NHDPlus Version 2: User Guide. 2012.
- Metherell, A. K., Harding, L. A., Cole, C. V., and Parton, W. J.: CENTURY Soil Organic Matter Model Environment.  
880 Technical Documentation Agroecosystem Version 4.0. Great Plains System Research Unit. Technical Report No. 4., Fort  
Collins, 1993.
- Nakhavali, M., Friedlingstein, P., Lauerwald, R., Tang, J., Chadburn, S., Camino-Serrano, M., Guenet, B., Harper, A.,  
Walmsley, D., Peichl, M., and Gielen, B.: Representation of dissolved organic carbon in the JULES land surface model (vn4.4-  
JULES-DOCM), *Geosci Model Dev*, 11, 593–609, <https://doi.org/10.5194/gmd-11-593-2018>, 2018.
- 885 Parton, W. J., Hartman, M., Ojima, D., and Schimel, D.: DAYCENT and its land surface submodel: description and testing,  
*Global and Planetary Change*, 35–48 pp., 1998.
- Parton, W. J., Schimel, D. S., Cole, C. V., and Ojima, D. S.: Analysis of Factors Controlling Soil Organic Matter Levels in  
Great Plains Grasslands, *Soil Science Society of America Journal*, 51, 1173–1179,  
<https://doi.org/10.2136/sssaj1987.03615995005100050015x>, 1987.
- 890 Poggio, L., De Sousa, L. M., Batjes, N. H., Heuvelink, G. B. M., Kempen, B., Ribeiro, E., and Rossiter, D.: SoilGrids 2.0:  
Producing soil information for the globe with quantified spatial uncertainty, *SOIL*, 7, 217–240, [https://doi.org/10.5194/soil-7-  
217-2021](https://doi.org/10.5194/soil-7-217-2021), 2021.
- Pool, S., Vis, M., and Seibert, J.: Evaluating model performance: towards a non-parametric variant of the Kling-Gupta  
efficiency, *Hydrological Sciences Journal*, 63, 1941–1953, <https://doi.org/10.1080/02626667.2018.1552002>, 2018.
- 895 Qualls, R. G. and Haines, B. L.: Biodegradability of Dissolved Organic Matter in Forest Throughfall, Soil Solution, and Stream  
Water, *Soil Science Society of America Journal*, 56, 578–586, <https://doi.org/10.2136/sssaj1992.03615995005600020038x>,  
1992.
- Romeiko, X. X., Guo, Z., Pang, Y., Lee, E. K., and Zhang, X.: Comparing machine learning approaches for predicting spatially  
explicit life cycle global warming and eutrophication impacts from corn production, *Sustainability (Switzerland)*, 12,  
900 <https://doi.org/10.3390/su12041481>, 2020.
- Ross, C. W., Prihodko, L., Anchang, J., Kumar, S., Ji, W., and Hanan, N. P.: HYSOGs250m, global gridded hydrologic soil  
groups for curve-number-based runoff modeling, *Sci Data*, 5, 180091, <https://doi.org/10.1038/sdata.2018.91>, 2018.
- Saltelli, A., Annoni, P., Azzini, I., Campolongo, F., Ratto, M., and Tarantola, S.: Variance based sensitivity analysis of model  
output. Design and estimator for the total sensitivity index, *Comput Phys Commun*, 181, 259–270,  
905 <https://doi.org/10.1016/j.cpc.2009.09.018>, 2010.
- Saltelli, A.: Making best use of model evaluations to compute sensitivity indices, *Comput Phys Commun*, 145, 280–297,  
[https://doi.org/10.1016/S0010-4655\(02\)00280-1](https://doi.org/10.1016/S0010-4655(02)00280-1), 2002.
- Santos, L., Thirel, G., and Perrin, C.: Technical note: Pitfalls in using log-transformed flows within the KGE criterion, *Hydrol  
Earth Syst Sci*, 22, 4583–4591, <https://doi.org/10.5194/hess-22-4583-2018>, 2018.
- 910 Scheller, J., Eklöf, K., Bishop, K., and Laudon, H.: Effects of forestry operations on dissolved organic carbon concentrations  
and export in boreal first-order streams, *J Geophys Res Biogeosci*, 117, <https://doi.org/10.1029/2011JG001827>, 2012.

- Sinsabaugh, R. L.: Phenol oxidase, peroxidase and organic matter dynamics of soil, <https://doi.org/10.1016/j.soilbio.2009.10.014>, March 2010.
- Sivapalan, M.: Pattern, Process and Function: Elements of a Unified Theory of Hydrology at the Catchment Scale, in: Encyclopedia of Hydrological Sciences, Wiley, <https://doi.org/10.1002/0470848944.hsa012>, 2005.
- 915 <https://doi.org/10.1002/0470848944.hsa012>, 2005.
- [Sobczak, W. V., Findlay, S., and Dye, S.: Relationships between DOC bioavailability and nitrate removal in an upland stream: An experimental approach, 2003.](https://doi.org/10.1002/0470848944.hsa012)
- Sobol, I. M.: Global sensitivity indices for nonlinear mathematical models and their Monte Carlo estimates, Mathematics and Computers in Simulation, 271–280 pp., 2001.
- 920 Steiner, J. L., Sadler, E. J., Chen, J.-S., Wilson, G., James, D., Vandenberg, B., Ross, J., Oster, T., and Cole, K.: Sustaining the ~~Earth's~~Earth's Watersheds-Agricultural Research Data System: Overview of development and challenges, J Soil Water Conserv, 63, 569–576, <https://doi.org/10.2489/jswc.63.6.569>, 2008.
- [Strauss, E. A. and Lamberti, G. A.: Effect of dissolved organic carbon quality on microbial decomposition and nitrification rates in stream sediments, Freshw Biol, 47, 65–74, https://doi.org/10.1046/j.1365-2427.2002.00776.x, 2002.](https://doi.org/10.2489/jswc.63.6.569)
- 925 Teodoru, C. R., Nyoni, F. C., Borges, A. V., Darchambeau, F., Nyambe, I., and Bouillon, S.: Dynamics of greenhouse gases (CO<sub>2</sub>, CH<sub>4</sub>, N<sub>2</sub>O) along the Zambezi River and major tributaries, and their importance in the riverine carbon budget, Biogeosciences, 12, 2431–2453, <https://doi.org/10.5194/bg-12-2431-2015>, 2015.
- Tian, H., Ren, W., Yang, J., Tao, B., Cai, W. J., Lohrenz, S. E., Hopkinson, C. S., Liu, M., Yang, Q., Lu, C., Zhang, B., Banger, K., Pan, S., He, R., and Xue, Z.: Climate extremes dominating seasonal and interannual variations in carbon export from the
- 930 Mississippi River Basin, Global Biogeochem Cycles, 29, 1333–1347, <https://doi.org/10.1002/2014GB005068>, ~~2015b~~2015a.
- Tian, H., Yang, Q., Najjar, R. G., Ren, W., Friedrichs, M. A. M., Hopkinson, C. S., and Pan, S.: Anthropogenic and climatic influences on carbon fluxes from eastern North America to the Atlantic Ocean: A process-based modeling study, J Geophys Res Biogeosci, 120, 752–772, <https://doi.org/10.1002/2014JG002760>, ~~2015a~~2015b.
- [Todd-Brown, K. E. O., Randerson, J. T., Post, W. M., Hoffman, F. M., Tarnocai, C., Schuur, E. A. G., and Allison, S. D.: Causes of variation in soil carbon simulations from CMIP5 Earth system models and comparison with observations, Biogeosciences, 10, 1717–1736, https://doi.org/10.5194/bg-10-1717-2013, 2013.](https://doi.org/10.1002/2014JG002760)
- 935 <https://doi.org/10.5194/bg-10-1717-2013>, 2013.
- Tranvik, L. J. and Jansson, M.: Terrestrial export of organic carbon, Nature, 415, 861–862, <https://doi.org/10.1038/415861b>, 2002.
- U.S. Geological Survey, National Water Information System data available on the World Wide Web (Water-Quality Data for the Nation): <https://waterdata.usgs.gov/nwis/qw>, <https://waterdata.usgs.gov/nwis/qw>, last access: 27 January 2024.
- USEPA, STorage and RETrieval Data Available on the World Wide Web (EPA STORET): <https://www.epa.gov/waterdata/storage-and-retrieval-and-water-quality-exchange>, <https://www.epa.gov/waterdata/storage-and-retrieval-and-water-quality-exchange>, last access: 27 January 2024.

945 (USGS), Environmental Protection Agency (EPA); 2021. <https://doi.org/10.5066/P9QRKUVJ>.

Wickland, K. P., Aiken, G. R., Butler, K., Dornblaser, M. M., Spencer, R. G. M., and Striegl, R. G.: Biodegradability of dissolved organic carbon in the Yukon River and its tributaries: Seasonality and importance of inorganic nitrogen, *Global Biogeochem Cycles*, 26, <https://doi.org/10.1029/2012GB004342>, 2012.

Wieczorek, M. E., Jackson, S. E., and Schwarz, G. E.: Select Attributes for NHDPlus Version 2.1 Reach Catchments and Modified Network Routed Upstream Watersheds for the Conterminous United States (ver. 3.0, January 2021): US Geological Survey data release, <https://doi.org/10.5066/F7765D7V>, 2018.

Wilson, H. F., Saiers, J. E., Raymond, P. A., and Sobczak, W. V.: Hydrologic Drivers and Seasonality of Dissolved Organic Carbon Concentration, Nitrogen Content, Bioavailability, and Export in a Forested New England Stream, *Ecosystems*, 16, 604–616, <https://doi.org/10.1007/s10021-013-9635-6>, 2013.

955 Yao, Y., Tian, H., Pan, S., Najjar, R. G., Friedrichs, M. A. M., Bian, Z., Li, H. Y., and Hofmann, E. E.: Riverine Carbon Cycling Over the Past Century in the Mid-Atlantic Region of the United States, *J Geophys Res Biogeosci*, 126, <https://doi.org/10.1029/2020JG005968>, 2021.

Ye, S., Li, H. Y., Huang, M., Alebachew, M. A., Leng, G., Leung, L. R., Wang, S. wen, and Sivapalan, M.: Regionalization of subsurface stormflow parameters of hydrologic models: Derivation from regional analysis of streamflow recession curves, *J Hydrol (Amst)*, 519, 670–682, <https://doi.org/10.1016/j.jhydrol.2014.07.017>, 2014.

960 Yeo, I. K. and Johnson, R. A.: A new family of power transformations to improve normality or symmetry, *Biometrika*, 87, 954–959, <https://doi.org/10.1093/biomet/87.4.954>, 2000.

Ying, X.: An Overview of Overfitting and its Solutions, in: *Journal of Physics: Conference Series*, <https://doi.org/10.1088/1742-6596/1168/2/022022>, 2019.

965 Zhang, P.: A novel feature selection method based on global sensitivity analysis with application in machine learning-based prediction model, *Applied Soft Computing Journal*, 85, <https://doi.org/10.1016/j.asoc.2019.105859>, 2019.

Zhao, M., Golaz, J. C., Held, I. M., Guo, H., Balaji, V., Benson, R., Chen, J. H., Chen, X., Donner, L. J., Dunne, J. P., Dunne, K., Durachta, J., Fan, S. M., Freidenreich, S. M., Garner, S. T., Ginoux, P., Harris, L. M., Horowitz, L. W., Krasting, J. P., Langenhorst, A. R., Liang, Z., Lin, P., Lin, S. J., Malyshev, S. L., Mason, E., Milly, P. C. D., Ming, Y., Naik, V., Paulot, F., Paynter, D., Philipps, P., Radhakrishnan, A., Ramaswamy, V., Robinson, T., Schwarzkopf, D., Seman, C. J., Shevliakova, E., Shen, Z., Shin, H., Silvers, L. G., Wilson, J. R., Winton, M., Wittenberg, A. T., Wyman, B., and Xiang, B.: The GFDL Global Atmosphere and Land Model AM4.0/LM4.0: 2. Model Description, Sensitivity Studies, and Tuning Strategies, *J Adv Model Earth Syst*, 10, 735–769, <https://doi.org/10.1002/2017MS001209>, 2018.

*Supplement of*

**~~Deriving a~~ Transformation Rate ~~Map~~Maps of Dissolved Organic Carbon ~~over~~in the Contiguous U.S.**

**Lingbo Li et al.**

5 *Correspondence to:* Hong-Yi Li ([hongyili.jadison@gmail.com](mailto:hongyili.jadison@gmail.com))

The copyright of individual parts of the supplement might differ from the article licence.

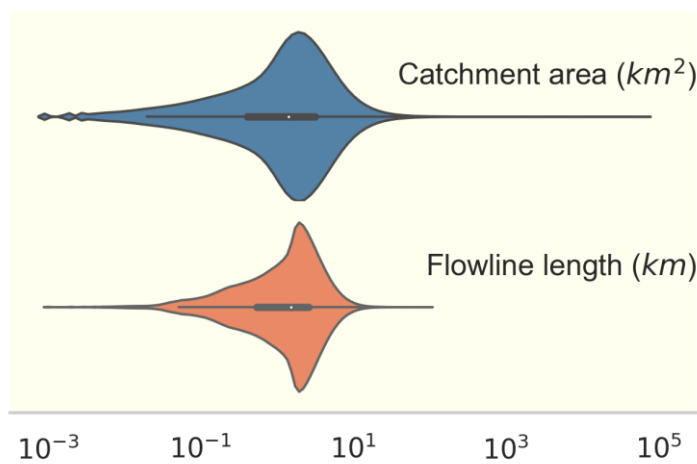


Figure S1. Distribution of the 2.6 million NHDPlus local catchment areareas and length-of-flowlines,flowline lengths.

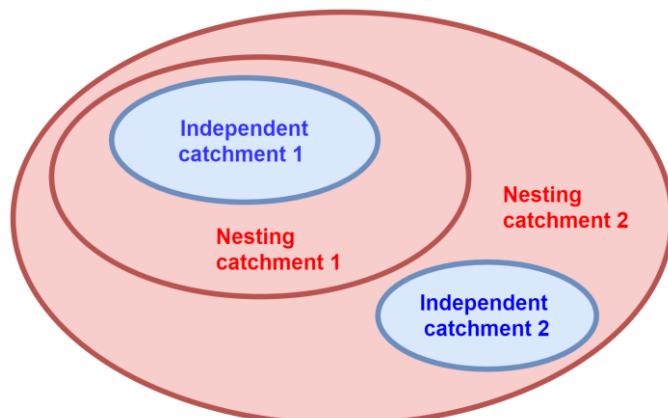
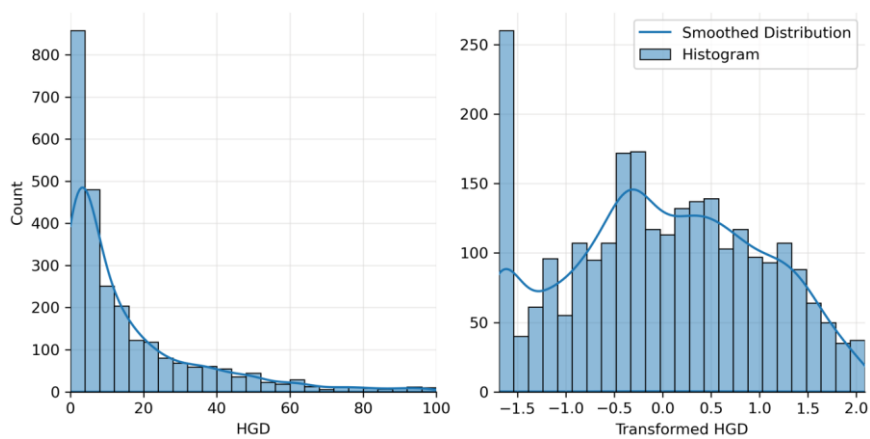


Figure S2. Comparison of catchment relationships: a) a small catchment and its containing NHDPlus local catchments, and b) independent catchment and nesting catchment. Note: The catchments outlined in red represent the same area, but only the boundary is shown in subplot a) for better visual clarity.

15



**Figure S3. Histogram of the percentage of hydrologic group D soil (HGD) predictor before (left panel) and after (right panel) power transformation.**



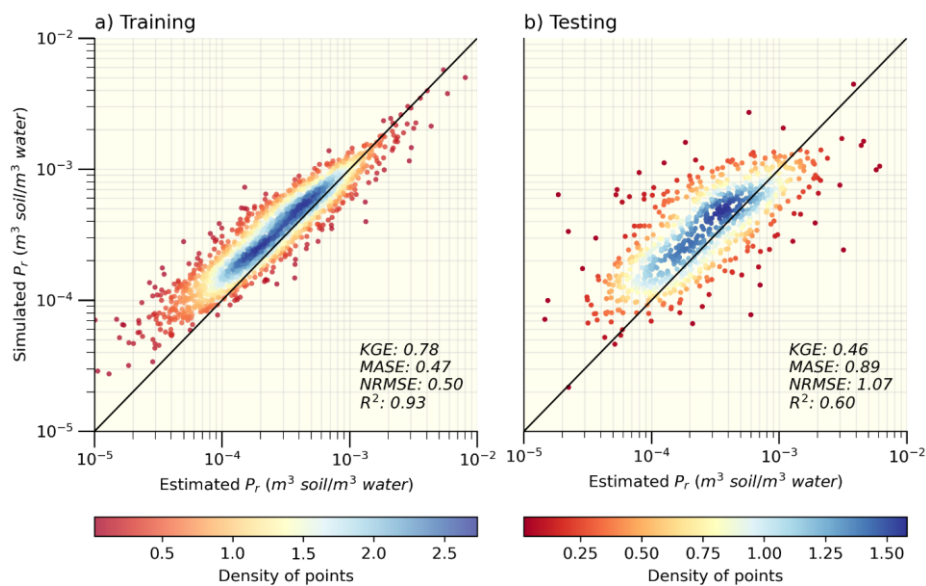
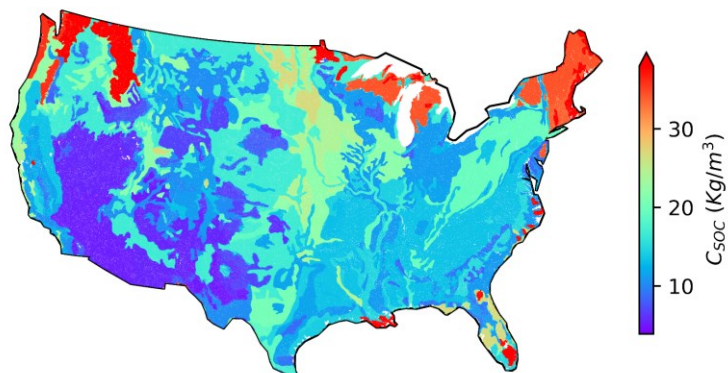
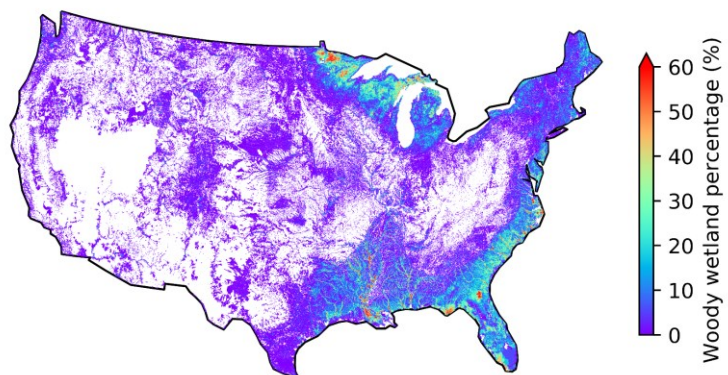


Figure S4. Performance of the XGBoost model (using HWSD) with 12 predictors trained by using KGE during a) the training phase (n=1816) and b) the testing phase (n=779). The solid black line indicates a 1:1 ratio. Note that the axes are in a log-log scale.

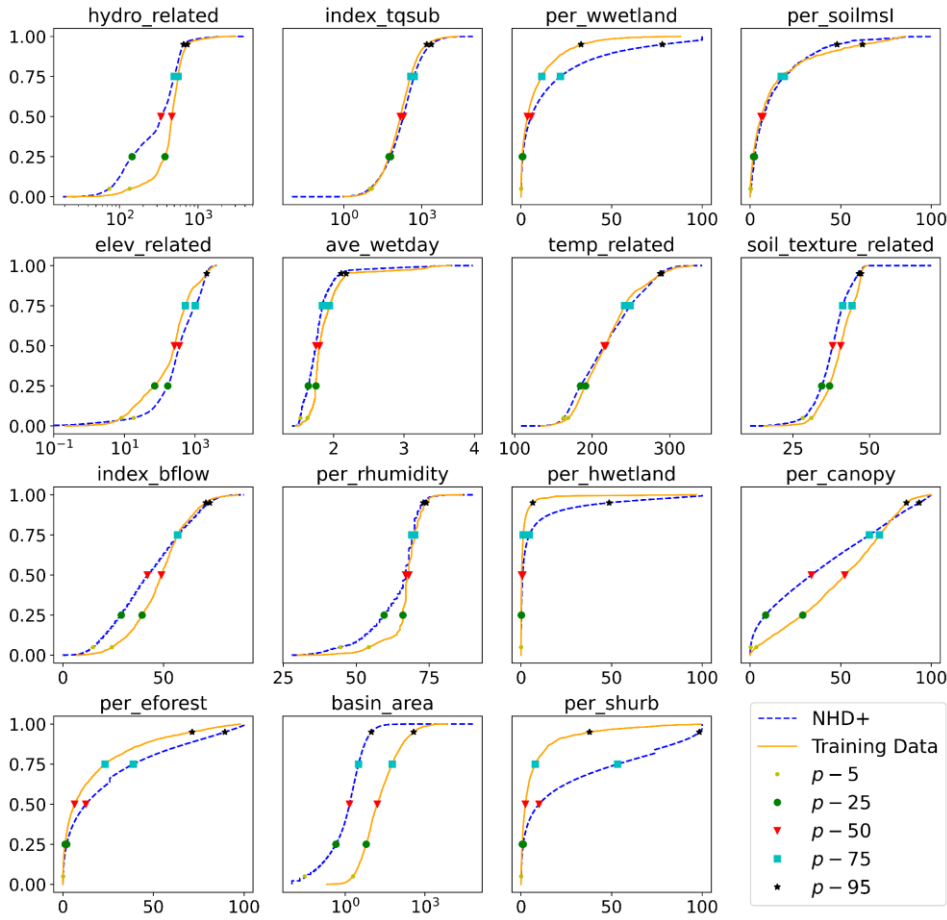
a)  $C_{SOC}$



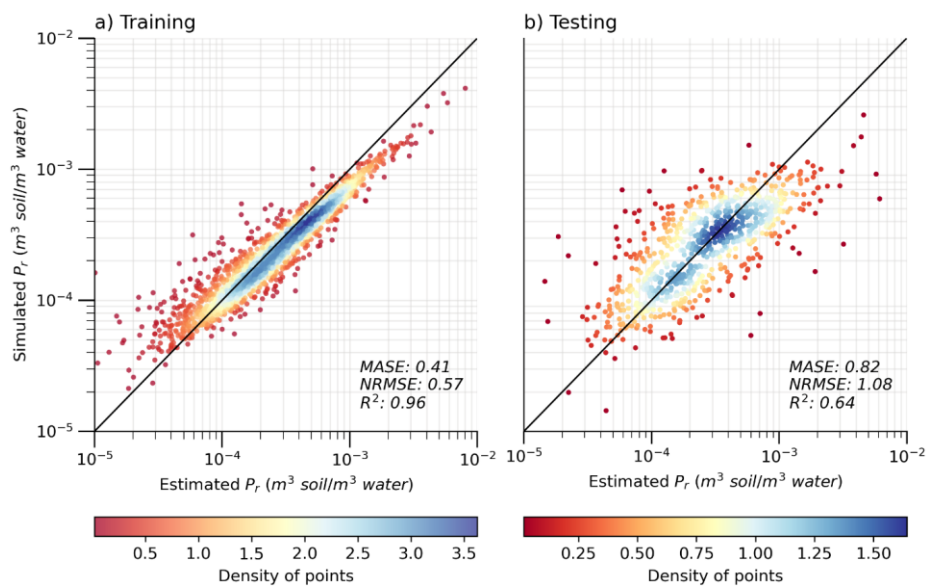
b) Woody wetland percentage



25 Figure S5. CONUS maps of: a) HWSO top-layer soil organic carbon (SOC) concentration, and b) woody wetland fraction across over 2 million NHDPlus catchments. Regions displayed in white may indicate missing data or a zero fraction.

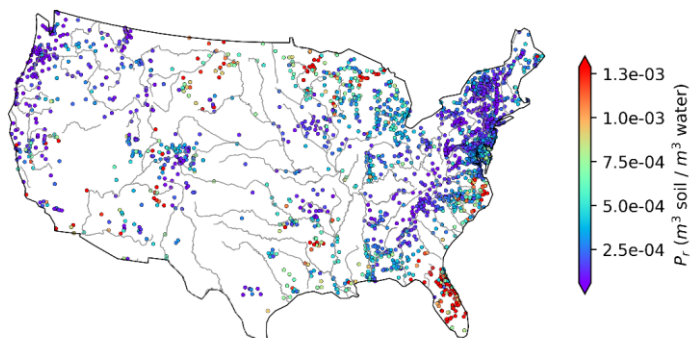


**Figure S6. Comparison of the cumulative distribution function (CDF) of 15 selected predictors between training data and all flowlines (i.e., NHDPlus). Abbreviations:** **hydro\_related** (merged predictor representing recharge, runoff, and precipitation); **CONTACT** (subsurface contact time); **NLCD01\_90** (areal percentage of woody wetlands); **HGBD** (areal percentage of Hydrologic Group BD soil); **elev\_related** (merged predictor for mean/min/max elevation); **CWD** (consecutive wet days); **temp\_related** (merged predictor encompassing potential evapotranspiration, first/last freeze timing, snow fraction, actual evapotranspiration, and mean/min/max temperature); **soil\_texture\_related** (merged predictor for silt and sand content); **BF1** (base flow index); **RH** (relative humidity); **NLCD01\_95** (areal percentage of herbaceous wetlands); **CNPY11\_BUFF100** (areal percentage of canopy in the riparian buffer); **NLCD01\_42** (areal percentage of evergreen forest); **BASIN\_AREA** (catchment area); **NLCD01\_52** (areal percentage of shrub). For detailed descriptions, refer to Supplementary Tables S1 and S2.



40 Figure S7. Performance of the XGBoost model (using HWSD) with 15 predictors during a) the training phase (n=1816) and b) the testing phase (n=779). The solid black line indicates a 1:1 ratio. Note that the axes are in a log-log scale.

a)  $P_r$  of independent catchments



b)  $P_r$  of evaluation catchments

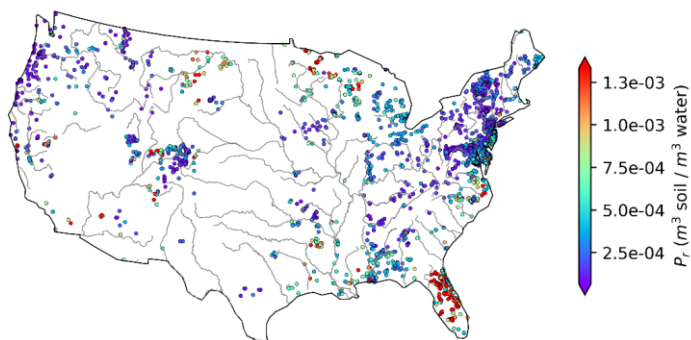
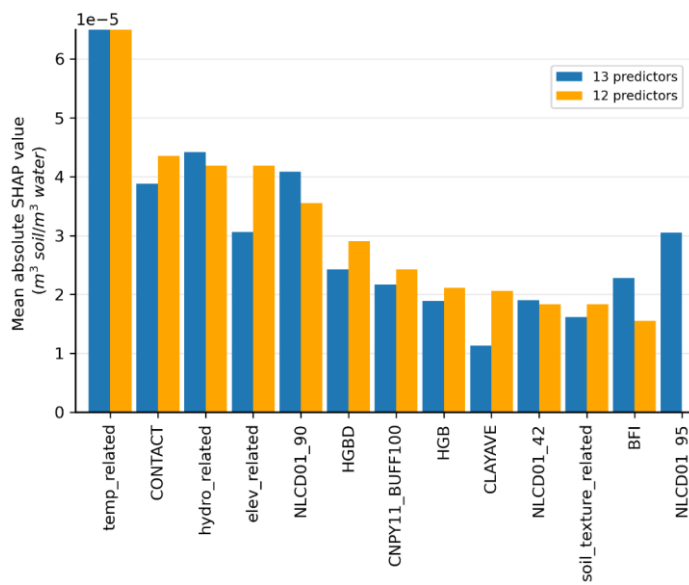
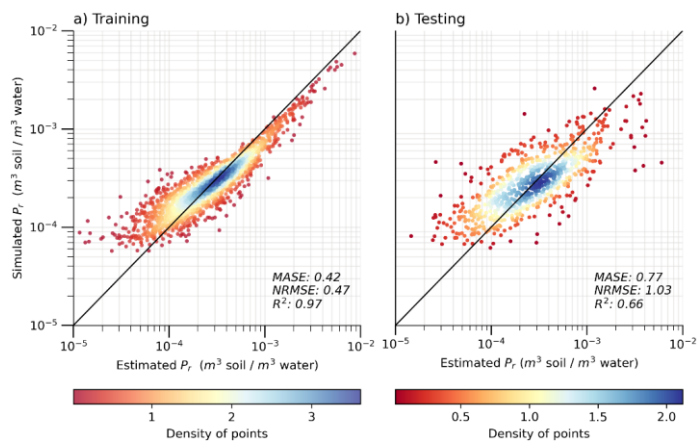


Table S1:

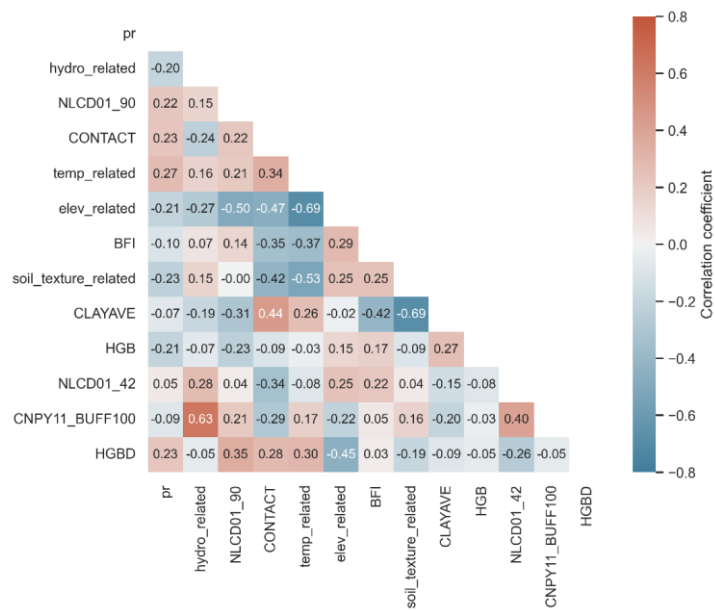
Figure S8. Variability in estimated  $P_r$  (using SoilGrids 2.0) across CONUS: a) For independent catchments ( $n=2583$ ), and b) For evaluation catchments ( $n=3210$ ). The points indicate the locations of the WQP stations, which are also the outlets of the corresponding small catchments. The CONUS boundary and river shapefiles are directly from open-source datasets GeoPandas (geopandas.org) and Natural Earth (Made with Natural Earth. Free vector and raster map data @ naturalearthdata.com), respectively. The color bars have been adjusted to enhance visual display by showing only the main body of values (from the 5th percentile to the 95th percentile).



**Figure S9. Mean absolute SHAP values of predictors in models (using SoilGrids 2.0) with 13 predictors (blue) and 12 predictors (orange).** Note that the SHAP values have the same units as the target variable,  $P_r$ . Abbreviations: temp\_related (merged predictor encompassing potential evapotranspiration, first/last freeze timing, snow fraction, actual evapotranspiration, and mean/min/max temperature); CONTACT (subsurface contact time); hydro\_related (merged predictor representing recharge, runoff, and precipitation); elev\_related (merged predictor for mean/min/max elevation); NLCD01\_90 (areal percentage of woody wetlands); HGBD (areal percentage of Hydrologic Group BD soil); CNPY11\_BUFF100 (areal percentage of canopy in the riparian buffer); HGB (areal percentage of Hydrologic Group B soil); CLAYAVE (clay content percentage); NLCD01\_42 (areal percentage of evergreen forest); soil\_texture\_related (merged predictor for silt and sand content); BFI (base flow index); NLCD01\_95 (areal percentage of herbaceous wetlands). For detailed descriptions, refer to Supplementary Tables S2 and S3.



**Figure S10. Performance of the XGBoost model (using SoilGrids 2.0) with 12 predictors during a) the training phase (n=1808) and b) the testing phase (n=775). The solid black line indicates a 1:1 ratio. The varying colours indicate the density of points in the scatter plot.**



**Figure S11. Covariance heatmap of  $P_r$  (using SoilGrids 2.0) and the 12 selected NHDPlus predictors. The Pearson correlation coefficient is used. Abbreviations: hydro\_related (merged predictor representing recharge, runoff, and precipitation); NLCD01\_90 (areal percentage of woody wetlands); CONTACT (subsurface contact time); temp\_related (merged predictor encompassing potential evapotranspiration, first/last freeze timing, snow fraction, actual evapotranspiration, and mean/min/max temperature); elev\_related (merged predictor for mean/min/max elevation); BFI (base flow index); soil\_texture\_related (merged predictor for silt and sand content); CLAYAVE (clay content percentage); HGB (areal percentage of Hydrologic Group B soil); NLCD01\_42 (areal percentage of evergreen forest); CNPY11\_BUFF100 (areal percentage of canopy in the riparian buffer); HGBD (areal percentage of Hydrologic Group BD soil). For detailed descriptions, refer to Supplementary Tables S2 and S3.**



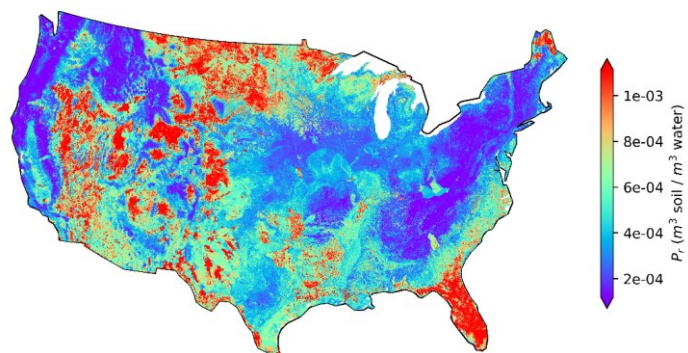
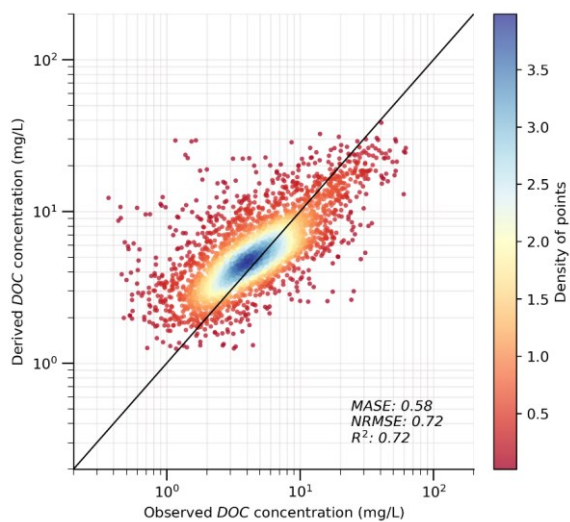
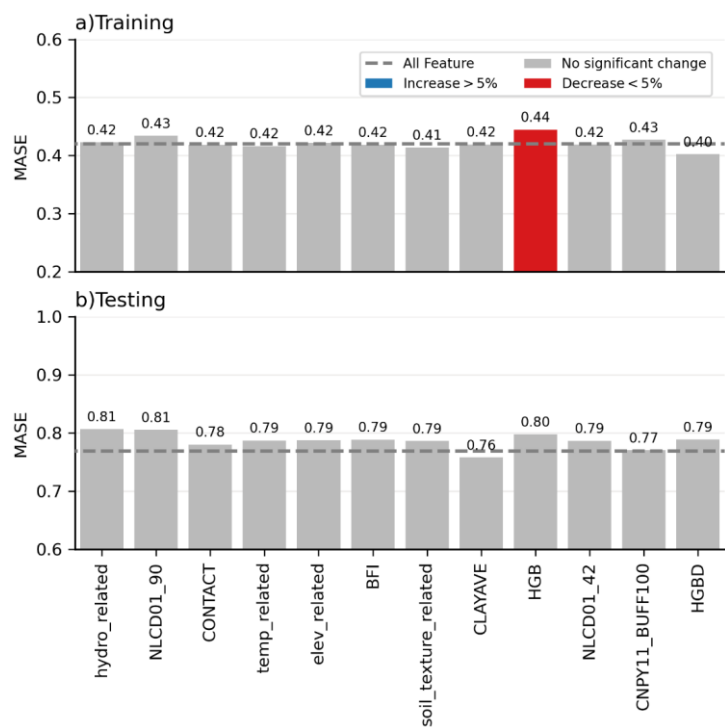


Figure S12. ML model (using SoilGrids 2.0) simulated  $P_r$  at over 2.6 million NHDPlus local catchments.



**Figure S13. Evaluation of derived DOC (using SoilGrids 2.0) concentration at the catchment scale (n=3208). The solid black line indicates a 1:1 ratio. The varving colours indicate the density of points in the scatter plot.**



**Figure S14. Sensitivity of XGBoost model (using SoilGrids 2.0) to predictors in the training and testing phases. The MASE value is represented by the blue, red, and grey bars, indicating whether the model performance increases, decreases, or remains relatively unchanged after dropping the corresponding predictor. The dashed grey line indicates the model performance with all variables included. Abbreviations: hydro\_related (merged predictor representing recharge, runoff, and precipitation); NLCD01\_90 (areal percentage of woody wetlands); CONTACT (subsurface contact time); temp\_related (merged predictor encompassing potential evapotranspiration, first/last freeze timing, snow fraction, actual evapotranspiration, and mean/min/max temperature); elev\_related (merged predictor for mean/min/max elevation); BFI (base flow index); soil\_texture\_related (merged predictor for silt and sand content); CLAYAVE (clay content percentage); HGB (areal percentage of Hydrologic Group B soil); NLCD01\_42 (areal percentage of evergreen forest); CNPY11\_BUFF100 (areal percentage of canopy in the riparian buffer); HGBD (areal percentage of Hydrologic Group BD soil). For detailed descriptions, refer to Supplementary Tables S2 and S3.**

**Table S1. In-stream DOC degradation rate (k) from previous modeling and experimental studies**

<u>Study Type</u>	<u>First-Order Decay rate (<math>k \text{ d}^{-1}</math>)</u>	<u>Study Domain</u>	<u>Reference</u>
<u>Modeling</u>	<u>0.01</u>	<u>Eastern North America</u>	<u>Tian et al., 2015</u>
	<u>0.01</u>	<u>Global</u>	<u>Li et al., 2019</u>
	<u>0.0163/0.0223<sup>a</sup></u>	<u>Upland and forested catchments in Canada</u>	<u>Futter et al., 2007</u>
<u>Experimental</u>	<u>0.011<sup>b</sup></u>	<u>Upland and forested catchment (Southern Appalachian Mountains, USA)</u>	<u>Qualls and Haines, 1992</u>
	<u>0.009<sup>b</sup></u>	<u>Upland and forested catchment (Catskill Mountains, USA)</u>	<u>Sobczak et al., 2003</u>
	<u>0.013<sup>c</sup></u>	<u>Forested headwater catchment (Hae-an Basin, South Korea)</u>	<u>Jung et al., 2014</u>
	<u>0.09<sup>c</sup></u>	<u>Agro-urban headwater catchments (Taihu Lake Watershed, China)</u>	<u>Wu et al., 2019</u>

*a. calibrated for the two catchments separately.*

*b. adopted from Table 2 in Mineau et al., 2016*

*c. calculated by fitting a first-order decay model using the published data.*

100 **Table S2.** List of 46 independent predictive attributes

Group	Acronym	Description
Chemical	NEMATOCIDE	Nematicide use on agricultural land (kg/km <sup>2</sup> )
Climate	CWD	Average number of consecutive days with measurable precipitation
	MINP6190	Watershed minimum average annual precipitation (mm)
	MINWD6190	Watershed average of minimum monthly number of days of measurable precipitation
	RH	Percent of the watershed average relative humidity
Geology	BEDPERM_1	Percent of NHDPlus version 2 flowline catchment whose bedrock permeability class is not a principal aquifer
	BEDPERM_6	Percent of NHDPlus version 2 flowline catchment whose bedrock permeability class is unconsolidated sand and gravel
	OLSON_PERM	Rock hydraulic conductivity (10 <sup>-6</sup> m/s)
	ROCKTYPE_200	Estimated percent of catchment that is underlain by the principal aquifer rock type, semi-consolidated sand aquifers
	ROCKTYPE_999	Estimated percent of catchment that is underlain by the principal aquifer rock type, other rocks
Hydrologic	BFI	Base Flow Index, a ratio of base flow to total streamflow, expressed as a percentage and ranging from 0 to 100
	CONTACT	The Subsurface flow contact time index estimates the duration infiltrated water resides in the saturated subsurface zone of the basin before discharging into the stream, measured in days
	IEOF	Percent of Horton overland flow
	SATOF	Percent of Dunne overland flow
Hydrologic Modifications	DITCHES92	Percent of watershed subjected to ditches for the year 1992
	MIRAD_2002	Percent of watershed in irrigated agriculture, from USGS 2002 250-m MODIS data

	MIRAD_2007	Percent of watershed in irrigated agriculture, from USGS 2007 250-m MODIS data
	MIRAD_2012	Percent of watershed in irrigated agriculture, from USGS 2012 250-m MODIS data
	TILES92	Percent of watershed subjected to tile drains for the year 1992
Land Cover	CNPY11_BUFF100	Percent of tree canopy in 100meter riparian buffer
	LAKEPOND	Percent of lakes or ponds
	NLCD01_11	Areal percent of 2001 land-use and land-cover type Open Water: All areas of open water, generally with less than 25 percent cover of vegetation or soil.
	NLCD01_31	Areal percent of 2001 land-use and land-cover type Barren Land: Barren areas of bedrock, desert pavement, scarps, talus, slides, volcanic material, glacial debris, sand dunes, strip mines, gravel pits, and other accumulations of earthen material. Generally, vegetation accounts for less than 15 percent of total cover.
	NLCD01_41	Areal percent of 2001 land-use and land-cover type Deciduous Forest: Areas dominated by trees generally greater than 5 meters tall, and greater than 20 percent of total vegetation cover. More than 75 percent of the tree species shed foliage simultaneously in response to seasonal change.
	NLCD01_42	Areal percent of 2001 land-use and land-cover type Evergreen Forest: Areas dominated by trees generally greater than 5 meters tall, and greater than 20 percent of total vegetation cover. More than 75 percent of the trees maintain their leaves all year. Canopy is never without green foliage.
	NLCD01_43	Areal percent of 2001 land-use and land-cover type Mixed Forest: Areas dominated by trees generally greater than 5 meters tall, and greater than 20 percent of total vegetation cover. Neither deciduous nor evergreen species are greater than 75 percent of total tree cover.
	NLCD01_52	Areal percent of 2001 land-use and land-cover type Shrub/Scrub: Areas dominated by shrubs less than 5 meters tall. Shrub canopy is typically greater than 20 percent of total vegetation. This class includes true shrubs, young trees in an early successional stage or trees stunted from environmental conditions.
	NLCD01_71	Areal percent of 2001 land-use and land-cover type Grassland/Herbaceous: Areas dominated by graminoid or herbaceous vegetation, generally greater than 80 percent of total vegetation. These areas are not subject to intensive management such as tilling, but can be utilized for grazing.

	NLCD01_81	Areal percent of 2001 land-use and land-cover type Pasture/Hay: Areas of grasses, legumes, or grass-legume mixtures planted for livestock grazing or the production of seed or hay crops, typically on a perennial cycle. Pasture/hay vegetation accounts for greater than 20 percent of total vegetation.
	NLCD01_90	Areal percent of 2001 land-use and land-cover type Woody Wetlands: Areas where forest or shrubland vegetation accounts for greater than 20 percent of vegetative cover and the soil or substrate is periodically saturated with or covered with water.
	NLCD01_95	Areal percent of 2001 land-use and land-cover type Emergent Herbaceous Wetlands: Areas where perennial herbaceous vegetation accounts for greater than 80 percent of vegetative cover and the soil or substrate is periodically saturated with or covered with water.
	SWAMPMARSH	Percent of swamps or marshes
Soils	CLAYAVE	Percent of average clay content
	HGA	Areal percent of Hydrologic Group A soil. Hydrologic group A soils have high infiltration rates. Soils are deep and well drained and, typically, have high sand and gravel content
	HGAD	Areal percent of Hydrologic Group AD soil. Hydrologic group AD soils have group A characteristics (high infiltration rates) when artificially drained and have group D characteristics (very slow infiltration rates) when not drained
	HGB	Areal percent of Hydrologic Group B soil. Hydrologic group B soils have moderate infiltration rates. Soils are moderately deep, moderately well drained, and moderately coarse in texture
	HGBD	Areal percent of Percentage of Hydrologic Group BD soil. Hydrologic group BD soils have group B characteristics (moderate infiltration rates) when artificially drained and have group D characteristics (very slow infiltration rates) when not drained
	HGC	Areal percent of Hydrologic Group C soil. Hydrologic group C soils have slow soil infiltration rates. The soil profiles include layers impeding downward movement of water and, typically, have moderately fine or fine texture
	HGCD	Areal percent of Hydrologic Group CD soil. Hydrologic group CD soils have group C characteristics (slow infiltration rates) when artificially drained and have group D characteristics (very slow infiltration rates) when not drained
	HGD	Areal percent of Percentage of Hydrologic Group D soil. Hydrologic group D soils have very slow infiltration rates. Soils are clayey, have a high water table, or have a shallow impervious layer

	SALINAVE	Salinity measured as average millimhos per centimeter
Topographic	ARTIFICIAL	Flowlines coded as Artificial paths
	BASIN_AREA	NHDPlusV2 flowline catchment area (km²)
	RDX	Number of road and stream intersections
	STREAMRIVER	Flowlines coded as Streams or Rivers
Water Use	FRESHWATER_WD	Freshwater withdrawals from 1995-2000 county-level estimates



Table 82S3. List of 49 correlated predictive attributes

Correlated Group	Acronym	Description
hydro_related	RECHG	Mean annual natural groundwater recharge (mm/yr).
	WB5100_ANN	Annual averaged runoff from McCabe and Wolock's Runoff Model 1951-2000 (mm).
	MAXP6190	Watershed maximum average annual precipitation (mm).
	PPT7100_ANN	Mean annual precipitation for the watershed, from 800m PRISM data (mm). 30 years period of record 1971-2000.
	RUN7100	Estimated 30-year (1971-2000) average annual runoff (mm/yr).
temp_related	PET	Mean annual potential evapotranspiration (PET), estimated using the Hamon equation.
	FSTFZ6190	Watershed average of mean day of the year of first freeze, derived from 30 years of record (1961-1990), 2km PRISM.
	LSTFZ6190	Watershed average of mean day of the year of last freeze, derived from 30 years of record (1961-1990), 2km PRISM.
	PRSNOW	Snow percent of total precipitation estimate, mean for period 1901-2000.
	ET	Mean-annual actual evapotranspiration (ET), estimated using regression equation of Sanford and Selnick (2013).
	TMAX7100	Watershed average of maximum monthly air temperature from 800m PRISM, derived from 30 years (1971-2000) of record (°C).
	TAV7100_ANN	Watershed average of monthly air temperature from 800m PRISM, derived from 30 years (1971-2000) of record (°C).
	TMIN7100	Watershed average of minimum monthly air temperature from 800m PRISM, derived from 30 years (1971-2000) of record (°C).
agri_chem_related	FUNGICIDE	Fungicide use on agricultural land (kg/km²).
	HERBICIDE	Herbicide use on agricultural land (kg/km²).

	INSECTICIDE	Insecticide use on agricultural land (kg/km <sup>2</sup> ).
	N97	Estimated nitrogen from fertilizer and manure.
	P97	Estimated phosphorous from fertilizer and manure.
	NLCD01_82	Areal percent of 2001 land-use and land-cover type Cultivated Crops: Areas used for the production of annual crops, such as corn, soybeans, vegetables, tobacco, and cotton, and also perennial woody crops such as orchards and vineyards. Crop vegetation accounts for greater than 20 percent of total vegetation. This class also includes all land being actively tilled.
	PEST219	Estimate of agricultural pesticide application (kg/km <sup>2</sup> ).
	KGBI	Toxicity Weighted Use for benthic invertebrates on agricultural land, 2013 (kg/km <sup>2</sup> ).
	KGCLADO	Toxicity Weighted Use for cladocerans on agricultural land, 2013 (kg/km <sup>2</sup> ).
	KGFISH	Toxicity Weighted Use for fish on agricultural land, 2013 (kg/km <sup>2</sup> ).
urban_related	POPDENS90	Population density from 1990 Census block level data, persons per km <sup>2</sup> per NHDPlus version 2 catchment.
	IMPV01_BUFF100	NLCD 2001 percent of imperviousness in the 100meter riparian buffer zones.
	IMPV06	NLCD 2006 percent of imperviousness per NHDPlus version 2 catchment.
	IMPV06_BUFF100	NLCD 2006 percent of imperviousness in the 100meter riparian buffer zones.
	POPDENS00	Population density from 2000 Census block level data, persons per km <sup>2</sup> per NHDPlus version 2 catchment.
	POPDENS10	Population density from 2010 Census block level data, persons per km <sup>2</sup> per NHDPlus version 2 catchment.

	NLCD01_21	Areal percent of 2001 land-use and land-cover type Developed, Open Space: Includes areas with a mixture of some constructed materials, but mostly vegetation in the form of lawn grasses. Impervious surfaces account for less than 20 percent of total cover. These areas most commonly include large-lot single-family housing units, parks, golf courses, and vegetation planted in developed settings for recreation, erosion control, or aesthetic purposes.
	NLCD01_22	Areal percent of 2001 land-use and land-cover type Developed, Low Intensity: Includes areas with a mixture of constructed materials and vegetation. Impervious surfaces account for 20-49 percent of total cover. These areas most commonly include single-family housing units.
	NLCD01_23	Areal percent of 2001 land-use and land-cover type Developed, Medium Intensity: Includes areas with a mixture of constructed materials and vegetation. Impervious surfaces account for 50-79 percent of the total cover. These areas most commonly include single-family housing units.
	NLCD01_24	Areal percent of 2001 land-use and land-cover type Developed, High Intensity: Includes highly developed areas where people reside or work in high numbers. Examples include apartment complexes, row houses, and commercial/industrial. Impervious surfaces account for 80-100 percent of the total cover.
	TOTAL_ROAD_DENS	Density of all road types per NHDPlusV2 catchment. Density is defined as the length of road divided by the catchment area.
	HDENS10	Historic housing densities for 2010.
soil_texture_related	SILTAVE	Percent of average value of silt content.
	SANDAVE	Percent of average value of sand content.
soil_restrictive_related	SRL25AG	Estimated percent of the soil restrictive layer in the upper 25cm of agricultural land.
	SRL35AG	Estimated percent of the soil restrictive layer in the upper 35cm of agricultural land.
	SRL45AG	Estimated percent of the soil restrictive layer in the upper 45cm of agricultural land.
	SRL55AG	Estimated percent of the soil restrictive layer in the upper 55cm of agricultural land.
wetd_related	MAXWD6190	Watershed average of maximum monthly number of days of measurable precipitation, derived from 30 years of record (1961-1990), 2.3km PRISM.

	WDANN	NHDPlusV2 catchment value for the 30year annual average (1961-1990) number of days of measurable precipitation.
topo_related	EWT	Average depth to water table relative to the land surface (m).
	TWI	Topographic wetness index, $\ln(a/S)$ ; where $\ln$ is the natural log, $a$ is the upslope area per unit contour length and $S$ is the slope at that point.
	BASIN_SLOPE	NHDPlusV2 flowline catchment's average slope in percent.
elev_related	ELEV_MEAN	NHDPlusV2 flowline catchment's mean elevation (m).
	ELEV_MIN	NHDPlusV2 flowline catchment's minimum elevation (m).
	ELEV_MAX	NHDPlusV2 flowline catchment's maximum elevation (m).

**Table S4. The optimal hyperparameters values of the XGBoost model (using SoilGrids 2.0).**

Hyperparameter	Optimal Value	Tuning Range	Default value	Description
<u>lambda</u>	<u><math>8.497 \times 10^{-1}</math></u>	<u><math>[0, \infty]</math></u>	<u>1</u>	<u>Control L1 and L2 regularization; the larger the value, the more conservative the model will be</u>
<u>alpha</u>	<u><math>2.198 \times 10^{-2}</math></u>	<u><math>[0, \infty]</math></u>	<u>0</u>	
<u>gamma</u>	<u><math>9.045 \times 10^{-2}</math></u>	<u><math>[0, \infty]</math></u>	<u>0</u>	<u>Govern the model learning process by changing the step size shrinkage and minimum loss reduction; the larger the value, the more conservative the model will be</u>
<u>eta</u>	<u><math>1.146 \times 10^{-1}</math></u>	<u><math>(0, 1]</math></u>	<u>0.3</u>	
<u>colsample_bytree</u>	<u><math>5.004 \times 10^{-1}</math></u>	<u><math>(0, 1]</math></u>	<u>1</u>	<u>Control the subsample ratio of columns and training instances; a proper set of those values will prevent the model from over-fitting</u>
<u>subsample</u>	<u><math>9.730 \times 10^{-1}</math></u>	<u><math>(0, 1]</math></u>	<u>1</u>	
<u>min_child_weight</u>	<u><math>3.123 \times 10^{-1}</math></u>	<u><math>[0, \infty]</math></u>	<u>1</u>	<u>Determine the growth of the tree</u>
<u>max_depth</u>	<u>8</u>	<u><math>[0, \infty]</math></u>	<u>6</u>	

110 **Table S5. Representativeness of XGBoost model (using SoilGrids 2.0) input predictors over CONUS.**

Attributes	Relative difference in percentiles between $P_{\text{available}}$ and $\text{whole\_conus}$ data					Average
	5th	25th	50th	75th	95th	
NLCD01_95	0.667	0.667	0.853	1.144	1.528	0.972
CNPY11_BUFF100	1.686	1.098	0.429	0.080	0.078	0.674
NLCD01_90	0.769	0.307	0.448	0.612	0.806	0.589
NLCD01_42	0.689	0.548	0.648	0.501	0.224	0.522
elev_related	0.736	0.783	0.312	0.618	0.008	0.491
hydro_related	0.586	0.899	0.317	0.109	0.106	0.403
HGBD	0.920	0.276	0.159	0.098	0.252	0.341
CONTACT	0.165	0.140	0.248	0.294	0.410	0.251
CLAYAVE	0.233	0.412	0.249	0.178	0.140	0.243
HGB	0.232	0.308	0.180	0.160	0.093	0.194
BFI	0.477	0.306	0.152	0.002	0.028	0.193
soil_texture_related	0.096	0.071	0.068	0.070	0.015	0.064
temp_related	0.036	0.034	0.008	0.029	0.005	0.022

Abbreviations: NLCD01\_95 (areal percentage of herbaceous wetlands); CNPY11\_BUFF100 (areal percentage of canopy in the riparian buffer); NLCD01\_90 (areal percentage of woody wetlands); NLCD01\_42 (areal percentage of evergreen forest); elev\_related (merged predictor for mean/min/max elevation); hydro\_related (merged predictor representing recharge, runoff, and precipitation); HGBD (areal percentage of Hydrologic Group BD soil); CONTACT (subsurface contact time); CLAYAVE (clay content percentage); HGB (areal percentage of Hydrologic Group B soil); BFI (base flow index); soil\_texture\_related (merged predictor for silt and sand content); temp\_related (merged predictor encompassing potential evapotranspiration, first/last freeze timing, snow fraction, actual evapotranspiration, and mean/min/max temperature). For detailed descriptions, refer to Supplementary Tables S2 and S3.

115

Table S6. Sobol sensitivity analysis results for the 12 selected predictors (using SoilGrids 2.0).

Predictors	Total Indices (ST)	First Order Indices (S1)	Difference ((ST-S1)/ST)
temp_related	0.573	0.316	0.448
hydro_related	0.192	0.124	0.351
CONTACT	0.118	0.044	0.624
NLCD01_42	0.115	0.009	0.918
elev_related	0.102	0.036	0.643
CNPY11_BUFF100	0.098	0.016	0.833
HGB	0.075	0.007	0.910
NLCD01_90	0.057	0.028	0.507
CLAYAVE	0.049	0.001	0.978
BFI	0.047	0.003	0.938
soil_texture_related	0.032	0.000	1.000
HGBD	0.005	0.002	0.499

Abbreviations: temp\_related (merged predictor encompassing potential evapotranspiration, first/last freeze timing, snow fraction, actual evapotranspiration, and mean/min/max temperature); hydro\_related (merged predictor representing recharge, runoff, and precipitation); CONTACT (subsurface contact time); NLCD01\_42 (areal percentage of evergreen forest); elev\_related (merged predictor for mean/min/max elevation); CNPY11\_BUFF100 (areal percentage of canopy in the riparian buffer); HGB (areal percentage of Hydrologic Group B soil); NLCD01\_90 (areal percentage of woody wetlands); CLAYAVE (clay content percentage); BFI (base flow index); soil\_texture\_related (merged predictor for silt and sand content); HGBD (areal percentage of Hydrologic Group BD soil). For detailed descriptions, refer to Supplementary Tables S2 and S3.

**Reference**

Futter, M. N., Butterfield, D., Cosby, B. J., Dillon, P. J., Wade, A. J., and Whitehead, P. G.: Modeling the mechanisms that control in-stream dissolved organic carbon dynamics in upland and forested catchments, *Water Resour Res*, 43, <https://doi.org/10.1029/2006WR004960>, 2007.

Jung, B. J., Lee, J. K., Kim, H., and Park, J. H.: Export, biodegradation, and disinfection byproduct formation of dissolved and particulate organic carbon in a forested headwater stream during extreme rainfall events, *Biogeosciences*, 11, 6119–6129, <https://doi.org/10.5194/bg-11-6119-2014>, 2014.

Li, M., Peng, C., Zhou, X., Yang, Y., Guo, Y., Shi, G., and Zhu, Q.: Modeling Global Riverine DOC Flux Dynamics From 1951 to 2015, *J Adv Model Earth Syst*, 11, 514–530, <https://doi.org/10.1029/2018MS001363>, 2019.

Mineau, M. M., Wollheim, W. M., Buffam, I., Findlay, S. E. G., Hall, R. O., Hotchkiss, E. R., Koenig, L. E., McDowell, W. H., and Parr, T. B.: Dissolved organic carbon uptake in streams: A review and assessment of reach-scale measurements, <https://doi.org/10.1002/2015JG003204>, 1 August 2016.

Qualls, R. G. and Haines, B. L.: Biodegradability of Dissolved Organic Matter in Forest Throughfall, Soil Solution, and Stream Water, *Soil Science Society of America Journal*, 56, 578–586, <https://doi.org/10.2136/sssaj1992.03615995005600020038x>, 1992.

Sobczak, W. V., Findlay, S., and Dye, S.: Relationships between DOC bioavailability and nitrate removal in an upland stream: An experimental approach, 2003.

Tian, H., Ren, W., Yang, J., Tao, B., Cai, W. J., Lohrenz, S. E., Hopkinson, C. S., Liu, M., Yang, Q., Lu, C., Zhang, B., Banger, K., Pan, S., He, R., and Xue, Z.: Climate extremes dominating seasonal and interannual variations in carbon export from the Mississippi River Basin, *Global Biogeochem Cycles*, 29, 1333–1347, <https://doi.org/10.1002/2014GB005068>, 2015.

Wu, Z., Wu, W., Lin, C., Zhou, S., and Xiong, J.: Deciphering the origins, composition and microbial fate of dissolved organic matter in agro-urban headwater streams, *Science of the Total Environment*, 659, 1484–1495, <https://doi.org/10.1016/j.scitotenv.2018.12.237>, 2019.