

# A machine learning approach to driver attribution of dissolved organic matter dynamics in two contrasting freshwater systems

Daniel Mercado-Bettín<sup>1</sup>, Ricardo Paíz<sup>2,3</sup>, Valerie McCarthy<sup>2</sup>, Eleanor Jennings<sup>3</sup>, Elvira de Eyto<sup>4</sup>, Angeles M. Gallegos<sup>5</sup>, Mary Dillane<sup>4</sup>, Juan C. Garcia<sup>5</sup>, José J. Rodríguez<sup>5</sup>, and Rafael Marcé<sup>1</sup>

<sup>1</sup>Centre for Advanced Studies of Blanes, Spanish National Research Council, Carrer Accés Cala Sant Francesc, 14, 17300, Blanes, Spain

<sup>2</sup>School of History and Geography, Dublin City University, D09 YT18, Dublin 9, Ireland

<sup>3</sup>Centre for Freshwater and Environmental Studies, Dundalk Institute of Technology, A91 K584 Dundalk, Co. Louth, Ireland

<sup>4</sup>Fisheries & Ecosystem Advisory Services, Marine Institute, F28 PF65 Newport, Co. Mayo, Ireland

<sup>5</sup>Ens d'Abastament d'Aigua Ter-Llobregat, Ctra. Aiguës, 6, 08440, Cardedeu, Spain

**Correspondence:** Daniel Mercado-Bettín (daniel.mercado@ceab.csic.es)

1 **Abstract.** Predicting water quality variables in lakes is critical for effective ecosystem management under climatic and human  
2 pressures. Dissolved organic matter (DOM) serves as an energy source for aquatic ecosystems and plays a key role in their bio-  
3 geochemical cycles. However, predicting DOM is challenging due to complex interactions between multiple potential drivers  
4 in the aquatic environment and its surrounding terrestrial landscape. This study establishes an open and scalable workflow to  
5 identify potential drivers and predict fluorescent DOM (fDOM) in the surface layer of lakes by exploring the use of supervised  
6 machine learning models, including random forest, ~~extreme gradient boosting~~, ~~light gradient boosting~~, ~~catboosting~~boosting  
7 methods, k-nearest neighbors, support vector regression and linear model. ~~It~~The workflow was validated in two contrasting  
8 systems: one natural lake in Ireland with a relatively undisturbed catchment, and one reservoir in Spain with a more human-  
9 influenced catchment. A total of 24 potential drivers were obtained from global reanalysis data, and lake and river process-based  
10 modelling. ~~Partial dependence and~~ SHapley Additive exPlanations (SHAP) ~~analyses~~ were conducted for the most influential  
11 drivers identified, with soil moisture, soil temperature, and Julian day being common to both study sites. The best prediction  
12 was ~~found~~obtained when using the CatBoost model (during hold-out testing period, Irish site: KGE > ~~0.69~~0.68,  $r^2$  > ~~0.54~~0.50;  
13 Spanish site: KGE > 0.66,  $r^2$  > 0.54). Interestingly, when only using drivers from globally accessible climate and soil reanalysis  
14 data plus Julian day, the prediction capacity was maintained at both sites, showcasing potential for scalability. Our findings  
15 highlight the complex interplay of environmental drivers and processes that govern DOM dynamics in lakes, and contribute to  
16 the modelling of carbon cycling in aquatic ecosystems.

## 17 1 Introduction

18 Lakes are an essential component of global biogeochemical cycles, sustain biodiversity, and provide critical ecosystem services,  
19 e.g., water supply, fishing and irrigation. However, their water quality is increasingly at risk due to climatic change and human  
20 pressures (Bhateria and Jain, 2016). A key water quality variable is dissolved organic matter (DOM), which influences light  
21 penetration, energy, oxygen dynamics and nutrient availability in any lake (Solomon et al., 2015). The dynamics of DOM

22 in lakes are driven by both external processes in the terrestrial environment and internal processes. Land cover, climate, and  
23 topography regulate (either increase or decrease) carbon production in the catchment and carbon inputs into the lake (Li et al.,  
24 2015). In the water body, the quantity and quality of DOM are controlled by physical and biogeochemical mechanisms such as  
25 photodegradation, microbial processing and mixing dynamics, but can also be impacted by water abstraction or dam regulation  
26 (Xenopoulos et al., 2021).

27 Increases in the concentrations of DOM can affect ecosystem stability and human water use (e.g., raw drinking water quality)  
28 by reducing oxygen levels, altering microbial communities and nutrient cycling (Lake et al., 2000). DOM is also a precursor  
29 to disinfection byproducts (DBPs) during drinking water treatment, substances which have negative human health implica-  
30 tions (Li et al., 2014). Understanding the dynamics of DOM in lakes is essential for water quality management, especially as  
31 climate-driven processes are expected to increasingly influence DOM in freshwater systems (Creed et al., 2018). Moreover, the  
32 occurrence of extreme events such as eutrophication, algal blooms and hypoxic events, for which levels of DOM play a key  
33 role, is also expected to increase (Gobler, 2020). Hence, predicting DOM in lake water can improve water quality mitigation  
34 protocols and support adaptive water use management strategies.

35 Predicting DOM dynamics remains a challenge as it results from complex interactions in the environment, including multiple  
36 biogeochemical processes (Weyhenmeyer and Karlsson, 2009). Modelling tools offer an approach to simulate DOM in lake  
37 water. Process-based models have traditionally been used to better understand lake water quality, including DOM dynamics  
38 (McCullough et al., 2018). However, they generally require a large number of model parameters and governing equations,  
39 i.e., extensive parameterisation, to represent these dynamics. On the other hand, machine learning (ML) models do not rely  
40 on parameter calibration but instead incorporate large amounts of driver variables and data. This functionality can leverage  
41 the increasing amount of data being collected through satellite imagery, high frequency monitoring, and global climate and  
42 environmental modelling initiatives (Müller et al., 2024; Toming et al., 2020; Asadollah et al., 2025).

43 ML models have emerged as potential tools for modelling complex environmental variables, including those related to  
44 hydrology (Nearing et al., 2021) and water quality (Hanson et al., 2020). They have been recently employed in environmental  
45 applications, showing good predictive capabilities due to their ability to handle high-dimensional data, and capture nonlinear  
46 relationships (Li et al., 2016), for a diversity of parameters in lakes such as chlorophyll-a (Chen et al., 2024), turbidity (Zhang  
47 et al., 2021), and nutrient concentrations (e.g., phosphorus) (Hanson et al., 2020), suggesting potential for predicting DOM  
48 (Herzprung et al., 2020).

49 This study introduces a workflow for predicting fluorescent DOM (fDOM) (a proxy for DOM) in lakes using a suite of  
50 supervised ML models driven by potential drivers either extracted from reanalysis data (climate and soil variables) or outputs  
51 from lake and catchment process-based models. The workflow was tested in two different study sites, one in Ireland and one  
52 in Spain, that represent contrasting settings for both the potential drivers and DOM dynamics. Model performance was first  
53 evaluated at each site using the most influential drivers to predict fDOM. Subsequently, a second simulation was performed  
54 using a subset of these drivers, limited to those sourced from reanalysis data, to evaluate the predictive capacity of the model in  
55 the context of higher scalability. ~~The key research questions guiding this study were~~ A primary intention of the study is to use  
56 ML models to produce a robust, accessible, and reproducible workflow, rather than to conduct an extensive model benchmark

57 comparison. The aims of this study are: (1) ~~What are to find~~ the most influential drivers ~~of fDOM predictions, and how does~~  
58 ~~their importance vary between~~ (features) in predicting fDOM, and analyse their importance in two contrasting sites?; (2)  
59 How to test how accurately can supervised ML ~~approaches models~~ predict lake fDOM driven by reanalysis-based data, and  
60 hydrologic and lake ~~modelling outputs?~~ (3) ~~How easily can the workflow be reproduced and scaled to other sites?~~ modeling  
61 outputs.

## 62 2 Materials and methods

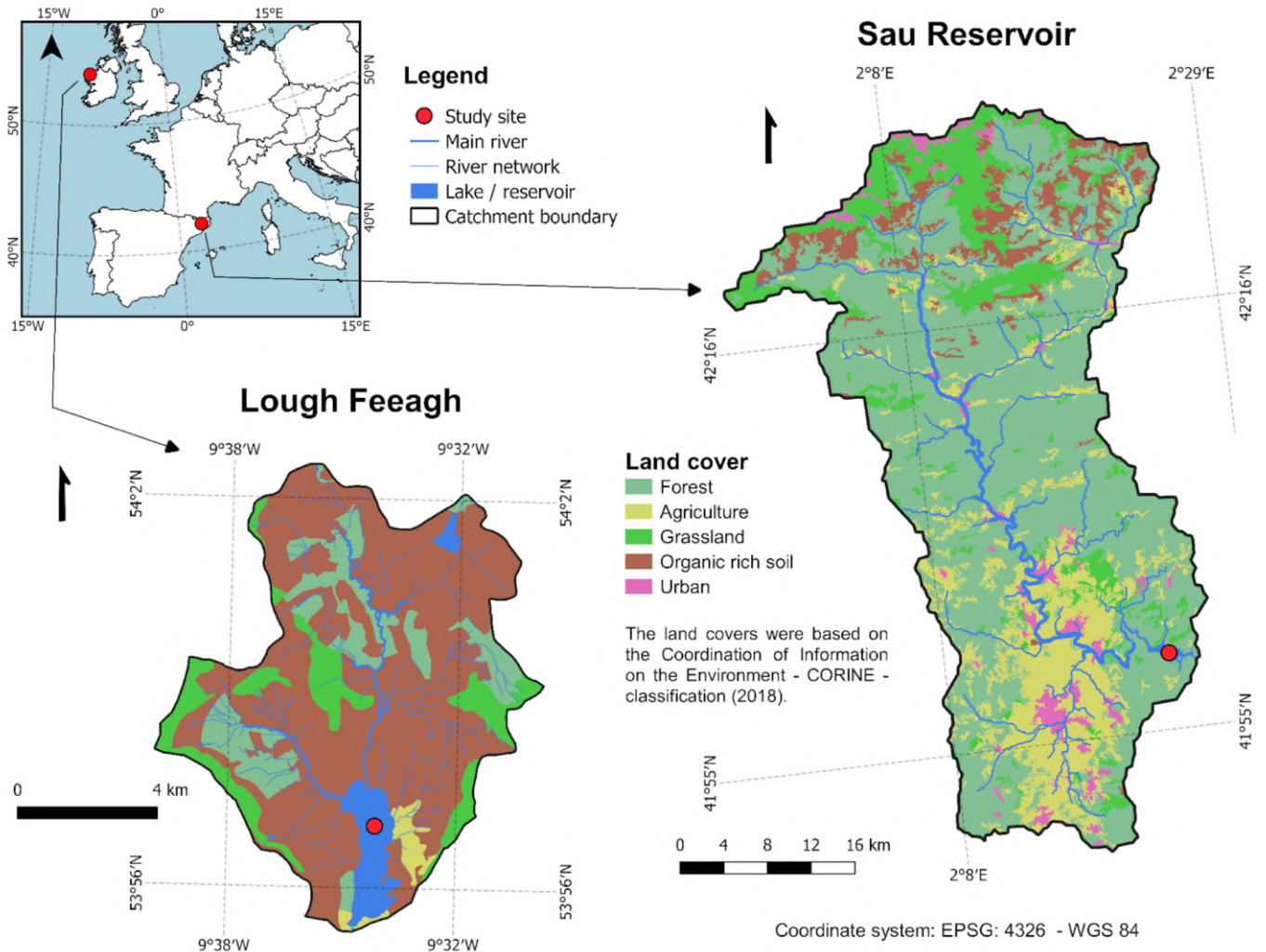
### 63 2.1 Study sites

64 Lough Feeagh and Sau Reservoir are located in western Ireland (53° 56' N, 9° 35' W) and northeastern Spain (41° 58' N, 2°  
65 23' E), respectively (Fig. 1). The study sites have contrasting attributes. Feeagh (depth of 46.8 m and area of 3.95 km<sup>2</sup>) is a  
66 monomictic and oligotrophic lake surrounded by a relatively undisturbed landscape, while Sau (depth of 70 m and area of 5.8  
67 km<sup>2</sup>) is an eutrophic system subjected to human activities and water abstraction. Feeagh has two primary inflows, the Black and  
68 the Glenamong rivers, while Sau has one, the Ter river. The catchment of Feeagh is relatively small (84 km<sup>2</sup>), with mid-range  
69 hills, and dominated by peatland. The catchment of Sau, in contrast, is larger (1525 km<sup>2</sup>), with a varying topography and land  
70 uses (Fig. 1; Fig. A1).

71 The dynamics of DOM in both study sites have been previously explored in Ryder et al. (2014) which identified that natural  
72 dynamics related to soil temperature, river discharge and drought ~~as were~~ important drivers in Feeagh, and in ~~Marcé et al.~~  
73 ~~(2021)~~ Marcé et al. (2021), which showed human activities ~~were significant drivers in Sau. DOM in Feeagh is mostly driven by~~  
74 ~~natural processes, while diffuse and point sources of organic matter and nutrients,~~ (e.g., wastewater effluents and agricultural  
75 runoff, ~~are also~~) were important for Sau, in addition to its environmental dynamics. Catchment hydrology is key for carbon  
76 transport into both study sites, and contributes to a distinct seasonality related to climate. Feeagh has a wet temperate climate,  
77 with cooler air temperatures and higher rainfall levels that occur on more than 75% of days in the year. The variability of carbon  
78 inputs reflects the sensitivity of a peatland-dominated landscape, which exacerbates climate-induced carbon release from the  
79 catchment into the lake. In contrast, Sau has a Mediterranean climate, characterised by hot, dry summers and mild, wet winters,  
80 dictating water availability, thermal stratification, and organic matter fluxes in the reservoir.

### 81 2.2 Prediction workflow

82 A five-step workflow was implemented to predict fDOM at each site (Fig. 2). First, all potential drivers for fDOM were  
83 collected as input data for each site. Second, an exploratory data analysis of fDOM was implemented and the ML models were  
84 trained using 85% of the available time series data (1978 out of 2328 fDOM measurements for Feeagh, and 653 out of 769 for  
85 Sau), while the remaining 15% (350 out of 2328 for Feeagh, and 116 out of 769 for Sau) future time series (hold-out period)  
86 was reserved for independent testing to evaluate performance and potential overfitting by comparing with the training period.



**Figure 1.** Two contrasting freshwater ecosystems— Lough Feeagh (Ireland) and Sau Reservoir (Spain) serve as contrasting sites for evaluating the predictive modelling of fDOM in lakes due to— The figure shows their distinct environmental locations and climatic conditions catchment land cover. The former Lough Feeagh is a humic and-oligotrophic lake — dominated by with a natural peatland peatland-dominated catchment — and temperate oceanic climate, resulting in relatively higher levels of DOC during the whole year with a regular seasonality; the latter whereas Sau is an- a eutrophic heavily controlled-reservoir — dominated by with a highly anthropized human-influenced catchment (urban wastewater effluents, intensive farming, agriculture), and a-Mediterranean climate— resulting in average in lower levels of DOC but with a much greater seasonal variation. Land cover data source: The CORINE Land Cover 2018 <https://doi.org/10.2909/960998c1-1870-4e82-8051-6485205ebbac> classes were grouped into the categories shown for visual clarity, following the aggregation approach described in the Supplementary Information (Figure A2).

87 Third, a set of drivers was selected for each site based on the variable importance extracted from the ML models, retaining  
88 only those drivers that exceeded an importance threshold of 5%. ~~A partial dependence and SHAP~~ SHAP and partial dependence  
89 analyses were applied to these specific drivers to evaluate how fDOM predictions at each site varied as a function of individually  
90 changing the selected input variables. Fourth, we ran simulations using (i) the drivers with a higher variable importance (> 5%),  
91 and (ii) using only drivers extracted from globally accessible reanalysis data, and assessed model ~~performance using coefficient~~  
92 ~~of determination (R<sup>2</sup>), Kling–Gupta efficiency (KGE), and root mean square error (RMSE).~~

93 The same workflow was applied to both sites using the same data sources, allowing for comparison. Following the FAIR  
94 principles, all data and workflow scripts are available and fully reproducible in the following repository: [https://github.com/  
95 danielmerbet/driver\\_attribution\\_fdom](https://github.com/danielmerbet/driver_attribution_fdom), a DOI will be provided in the final version of the manuscript. Large language models  
96 were used in this study to optimise the codes, improve the final plots and, for basic proofreading of the text.

## 97 **2.3 Data**

### 98 **2.3.1 Target variable (fDOM)**

99 Daily surface fDOM values were ~~computed~~ obtained from high-frequency data (~~2 minutes resolution~~) for both sites, ~~for~~  
100 ~~Feeagh~~ (2-minute-resolution data averaged to obtain daily values before analysis). For Feeagh, data were measured at 0.9 m  
101 depth, and for Sau, an average value was calculated between the ~~depths of measurements for depths~~ 0-5 m. fDOM data were  
102 expressed as quinine sulfate units (QSU) for the analysis. In Feeagh, the data spanned from 1st of May 2012 to 19th of  
103 November 2019 (n = 2328), and for Sau from 4th of February 2017 to 2nd of March 2020 (n = 769), with some gaps. All the  
104 other data (i.e., driver data) used in the workflow of this study were constrained by the availability of fDOM data. Following  
105 exploratory data analysis (Fig. A3), no transformation was applied to the fDOM data prior to training the ML models in the  
106 main manuscript, results for Sau are provided with a log<sub>10</sub> transformation in the Supplementary Information for comparison  
107 (this is revisited in Results), given the moderate skewness of fDOM at this study site.

108 In Feeagh, fDOM data were collected using a Seapoint UV fluorometer sensor (Seapoint Sensors Inc., Exeter, NH, USA)  
109 and water temperature data were measured using a Hach Environmental Hydrolab Data Sonde X5 (UK OTT Hydrometry Ltd).  
110 In Sau, fDOM and water temperature data were collected using a fDOM Digital Smart Sensor and Multiparameter Sonde  
111 (YSI EXO sonde, Yellow Springs, OH, USA), respectively. Raw fDOM data were water temperature-corrected in both sites  
112 based on relationships established for each sensor (Ryder et al., 2012). Details about the fDOM corrections can be found in  
113 Supplementary Information ~~and Figures A2 and A3~~ (Fig. A4, A5 and A6).

### 114 **2.3.2 Driver data**

115 ~~The input data for the workflow comprised 24 driving variables at each site.~~ All input data variables, including their respective  
116 units and source, are displayed in Step 1 of Figure 2. These were grouped into five categories: (1) meteorology, (2) soil, (3)  
117 process-based hydrological modelling, (4) process-based lake modelling and (5) Julian day. All input data variables, including  
118 their respective units and source, are displayed in Step 1 of Figure 2. Daily values of meteorology and soil variables were

119 extracted from the European Centre for Medium-Range Weather Forecasts (ECMWF) Reanalysis v5 (ERA5) dataset (Hersbach  
120 et al., Accessed on July 2025). This gridded dataset provides (pseudo) observations at a global scale with a spatial resolution  
121 of 0.25°. Eight meteorological variables (traditionally employed in water modelling studies) and soil temperature and soil  
122 moisture data at four depths (0-7 cm, 7-28 cm, 28-100 cm, and 100-255 cm) were extracted for the grid-cell which contained  
123 each water body.

124 Daily values of inflow discharge and inflow DOC concentration into each site were generated using the Generalised Water-  
125 shed Loading Functions Model (GWLf) coupled with a DOC module (GWLf-DOC). The GWLf-DOC version and calibration  
126 strategy applied are described in Paíz et al. (2025a). Calibration results can be found in the Supplementary Information. Daily  
127 values of five key lake variables (see Step 1, Fig. 2) were obtained from the General Lake Model (GLM) (Hipsey et al., 2019)  
128 run for both sites. Calibration strategies applied are described in Mercado-Bettín et al. (2021); Paíz et al. (2025b). Calibration  
129 results can be found in the Appendix. In addition, the cosine (to avoid an abrupt numerical change at every start of a year)  
130 of Julian day was included in the driver data inputs, given that seasonality is expected to influence DOM predictions; the  
131 approximate calendar timing corresponds to the cosine of the Julian day: values close to 1 correspond to boreal winter, 0 to  
132 spring/autumn, and -1 to boreal summer.

## 133 2.4 Supervised Machine Learning

134 Supervised ML models have advantages and limitations for time series prediction. In addition to capturing non-linear rela-  
135 tionships typical in aquatic systems and water quality predictions (Hollister et al., 2016; Regier et al., 2023), they provide  
136 flexibility to assess multiple drivers, temporal indicators, and variables external to the system (Qi, 2012; Rodriguez-Galiano  
137 et al., 2015). ML models do not require a fixed set of drivers to predict fDOM effectively, unlike process-based models,  
138 which typically rely on predefined inputs. Additionally, there is no need for parameter calibration but hyperparameter tun-  
139 ing, simplifying the modelling process. While some may argue that the lack of parameterization suggests a "black box" ap-  
140 proach, supervised ML can provide insights into the potential drivers for predicting a target variable (~~Biau and Scornet, 2016~~)  
141 (Biau and Scornet, 2016; Molnar, 2020).

142 However, due to the intrinsic autocorrelation in time series, e.g., when predicting DOM, these models tend to overfit when  
143 using out-of-bag samples during training. To overcome this issue, robust validation and training are required. Here, we used  
144 a hold-out period for validation during testing at both study sites. Prior to selecting this method, we compared it with two  
145 alternative validation methods using random forest: (1) 5-fold cross-validation and (2) rolling window cross-validation with  
146 a two-year training period, a one-year testing period, and a window shift every 90 days (see supplementary Figure A4). ~~The~~  
147 ~~comparison revealed that Feeagh exhibited more consistent model performance, with less overfitting between training and~~  
148 ~~testing phases, across the different validation methods, compared to Sau. This difference is likely attributable to the limited~~  
149 ~~amount of available data at Sau. A7).~~

150 ~~Seven ML and statistical models were used~~ Supervised machine-learning models were applied to predict fDOM dynamics:  
151 , including Random Forest (RF) (Breiman, 2001), eXtreme Gradient Boosting (XGBoost) (Chen and Guestrin, 2016), Light  
152 Gradient Boosting (LGB) (Ke et al., 2017), CatBoost (CTB) (Prokhorenkova et al., 2019), k-Nearest Neighbors (KNN) (Fix,

153 1985), Support Vector Regression (SVR) (Cortes and Vapnik, 1995) and linear model. ~~RF, XGBoost, LGB and CTB can~~  
154 ~~directly provide the importance of~~ The three gradient-boosting frameworks were included to assess the robustness of results  
155 across closely related implementations. Given their similarities, the main results shown are focused on the highest-performing  
156 model among the three. Neural networks were not included because they typically require more complex architecture design  
157 and calibration, which can increase methodological complexity and may reduce reproducibility of the ~~features to predict the~~  
158 ~~target variable, hence, only these four models were used to select the most important drivers to predict fDOM. For this, the~~  
159 ~~increase in node purity was used for RF, and the gain contribution of each feature to the model for XGB, LGB and CTB.~~  
160 ~~Further, hyperparameter tuning workflow.~~ In contrast, the selected models allow more reproducible implementation due to  
161 their comparatively simpler hyperparameter tuning.

#### 162 2.4.1 Hyperparameter tuning

163 Hyperparameter tuning was implemented in all ML models to improve accuracy and ~~generalisation. To extract the importance~~  
164 ~~of the drivers, implement hyperparameter tuning~~ generalization, and reduce the risk of overfitting. It was conducted in R using  
165 the caret framework where possible: for RF  $mtry \in \{2, 4, 6\}$ ; for XGBoost  $nrounds \in \{1000, 2000\}$ ,  $max\_depth \in \{3, 6\}$ ,  
166  $eta \in \{0.01, 0.05, 0.1\}$ ,  $gamma = 0$ ,  $colsample\_bytree = 1$ ,  $min\_child\_weight = 1$ , and  $subsample = 1$ ; for kNN  $k \in \{3, 5, 7\}$ ;  
167 for SVR  $C \in \{0.1, 1, 10\}$ , and  $sigma \in \{0.01, 0.1\}$ . Additional model-specific implementations were made: for LGB  $learning\_rate \in \{0.$   
168  $num\_leaves \in \{15, 31, 63\}$ ,  $max\_depth \in \{5, 10, 1\}$ ,  $feature\_fraction \in \{0.8, 1\}$ ,  $bagging\_fraction \in \{0.8, 1\}$ ,  $nrounds = 2000$   
169 (fixed); and for CTB  $iterations \in \{300, 500, 1000\}$ ,  $depth \in \{4, 6, 8\}$ ,  $learning\_rate \in \{0.01, 0.05, 0.1\}$ , and  $l2\_leaf\_reg \in \{1, 3, 5\}$ .  
170 Models were trained on the training subset and evaluated on the hold-out test set using RMSE. The parameter configuration  
171 minimising test RMSE was selected. The hyperparameter tuning can be reproduced by using the code “2\_hyperparameter\_tuning.R”  
172 in the repository.

173 For each case study, the dataset was chronologically split into training (85%) and hold-out test (15%) subsets. The final  
174 15% of observations were reserved as an independent test set to evaluate predictive performance, mimicking a forecasting  
175 exercise. Hyperparameter tuning was conducted exclusively on the training subset. Within the training data, model selection  
176 was performed using 5-fold cross-validation. The data were randomly partitioned into five equally sized folds; four folds were  
177 used for model fitting and one for validation, iteratively. The average cross-validated performance was used to select the optimal  
178 hyperparameter configuration.

#### 179 2.4.2 Performance metrics

180 To evaluate and facilitate inter-model comparison, three complementary performance metrics were computed for both training  
181 and test datasets: Coefficient of determination ( $R^2$ ), Root Mean Square Error (RMSE), and Kling–Gupta Efficiency (KGE).  
182 These metrics were selected because they capture complementary aspects of model performance such as correlation strength  
183 ( $R^2$ ), error magnitude (RMSE), and combined accuracy in correlation, bias, and ~~predict fDOM, multiple R packages were~~  
184 ~~used: caret, randomForest, xgboost, lightgbm, catboost, kknn and kernlab~~ variability (KGE). In addition, they are widely used  
185 in water-quality modelling of aquatic ecosystems, allowing comparison of the results with other studies.

## 186 2.5 ~~Partial-dependence-plots~~ Driver importance and SHAP analysis

187 To assess the influence of the most important drivers, we included ~~partial-dependence-and~~ SHapley Additive exPlanations  
188 (SHAP) plots, using the ~~Random-Forest-and-CatBoost-models,-respectively.~~ The partial-dependence-plots illustrate the individual  
189 influence-of-each-driver-on-fDOM-predictions-by-varying-the-driver's-values-across-its-entire-range-while-keeping-all-other  
190 drivers-constant-in-their-average-value ~~best boosting method~~. The SHAP plots measures how much a single driver (feature)  
191 value contributes to moving the prediction away from the average value, the Y-axis has the input drivers ranked by overall  
192 importance (from top to bottom), X-axis has the SHAP value representing the impact on model output for a single prediction,  
193 each point is a single data instance and the color reflects the driver value (blue = low, red = high). ~~To implement~~ In addition,  
194 in-supplementary-information, partial dependence plots ~~and SHAP-analysis,-the-pdp-package-in-R-and-the-were-generated~~  
195 to-support-and-illustrate-the-individual-influence-of-each-driver-on-fDOM-predictions-by-varying-the-driver's-values-across  
196 its-entire-range-while-keeping-all-other-drivers-constant-in-their-average-value. These partial dependence plots were generated  
197 using-the-outputs-from-the-Random-Forest-model. ~~To implement SHAP analysis and partial dependence plots, the shap package~~  
198 ~~in Python-were-used~~ Python and the pdp package in R were used, respectively.

## 199 3 Results

### 200 3.1 Driver attribution

201 The most influential predictors of fDOM at each site were identified from all 24 potential drivers based on the 5% threshold of  
202 the variable importance extracted from the ~~RF, XGB, LGB and CTB models~~ CTB model. This resulted in ~~eight~~ seven influential  
203 drivers being identified for Feeagh and five for Sau (Figure 3), four of which were common to both sites. This selection was  
204 consistent-across-multiple-machine-learning-models (Fig. A8).

205 The variables for the deepest soil layer were relevant for both study sites. Soil temperature and soil moisture at 100-255  
206 cm, were the most influential drivers for Feeagh. Similarly, for Sau, ~~the-deepest-soil-moisture-driver-was-remarkably-soil~~  
207 moisture-and-temperature-at-similar-depths-were the most influential ~~,-while-the-deepest-soil-temperature-was-still-important~~  
208 ~~but-less-so-than-in-Feeagh~~ drivers. Another key driver that was shared between Feeagh and Sau was Julian day. Lake volume  
209 was only important for Sau, while solar radiation, the amount of carbon entering the water body (indicated by the DOC inflow  
210 concentration) and both soil moisture and temperature at 28-100 cm were only influential for Feeagh.

### 211 3.2 ~~Partial-dependency~~ Driver influence on ~~selected-drivers~~ fDOM predictions

212 Figure 4 introduces ~~partial-dependence-plots-and~~ SHAP beeswarm plots for the most influential drivers selected in Figure 3,  
213 enabling the assessment of the individual effect of each driver on fDOM predictions.

214 Seasonal patterns in fDOM concentrations were observed in both Feeagh and Sau, with higher values in winter and lower  
215 in summer, as reflected in the influence of Julian day. However, the key predictors and their effects differed substantially,  
216 shaped by contrasting catchment and climate characteristics. In Feeagh, where precipitation is relatively high and sustained

217 year-round, deep soil temperature (100–255 cm) was the dominant and potentially limiting predictor, with fDOM increasing  
218 with temperature up to a threshold, beyond which a drop in the water table may counteract the effect ([see partial dependence  
219 plots; Fig. A9 and A10](#)). In addition, Feeagh showed minimal influence of mid-depth soil temperature (28-100 cm) and solar  
220 radiation on fDOM. This temperature-fDOM relationship was less relevant in Sau, where deep soil temperature (100-255 cm)  
221 had less explanatory power. Instead, fDOM dynamics in Sau, were driven primarily by soil moisture at both 28-100 cm and  
222 100-255 cm depths, potentially depicting a limiting condition by water stress.

223 Water availability also shaped the role of other predictors differently across the two sites. For instance, surface water tem-  
224 perature in Feeagh showed a clear threshold behaviour, with fDOM increasing relatively linearly beyond 6.5°C and stabilizing  
225 around 7.5°C, while in Sau, water volume acted as a surrogate for fDOM production ([see partial dependence plots; Fig. A9  
226 and A10](#)). Lower volumes corresponded to reduced fDOM values, which increase and stabilize beyond 146 hm<sup>3</sup>. Interest-  
227 ingly, inflow DOC concentration influences Feeagh more than Sau, likely due to differences in hydroclimatological processes  
228 governing these relationships. These differences highlight how catchment water availability fundamentally alters the relative  
229 importance and behavior of fDOM drivers.

### 230 3.3 Predicting DOM using Supervised Machine Learning

231 Table 1 presents a comparative evaluation of the seven ML and statistical models employed in this study. In Lough Feeagh,  
232 CatBoost (CTB) demonstrated the best overall performance with the highest R<sup>2</sup> ([0.510.50](#)) and KGE ([0.690.68](#)), and a relatively  
233 low RMSE ([8.298.17 QSU](#)). Similarly, in Sau Reservoir, CTB again has the highest R<sup>2</sup> (0.54) and KGE (0.66), and one of the  
234 lowest RMSE (7.11 [QSU](#)). While some models like RF and LM showed moderate performance at Feeagh, others such as SVR  
235 and XGB performed poorly in Sau, with SVR even having a negative KGE (-0.82), suggesting significant model bias. Overall,  
236 CatBoost (CTB) consistently outperformed all other models across both sites, supporting its selection for fDOM prediction  
237 using the selected environmental drivers.

238 The results during the training phase (supplementary Table [A1A2](#)) confirm that most ML models, especially XGB, RF,  
239 and KNN, had a high performance during training. However, the performance dropped in some models (e.g., XGB at Sau  
240 Reservoir) during the testing phases, underscoring the importance of evaluating models on independent test data to assess  
241 generalisability and overfitting. CatBoost (CTB), again, presented a more stable performance, showing slightly lower train-  
242 ing metrics (especially in Feeagh) compared with the other ML models and better generalisation when comparing with the  
243 metrics during testing, ~~supporting its suitability for prediction.~~ [Feeagh exhibited more consistent model performance, with  
244 less overfitting between training and testing phases, across the different validation methods \(Fig A7\), compared to Sau. This  
245 difference is likely attributable to the limited amount of available data at Sau. In addition, we assess whether this overfitting  
246 could be related to the moderate skewness of the fDOM data in Sau, however similar results and performance metrics were  
247 found when applying log transformation \(to remove skewness\) before applying the ML models \(see Fig. A11 and Tables A3  
248 and A4\).](#)

**Table 1.** Model comparison to predict fDOM using all selected drivers in both study sites. Statistic metrics ( $R^2$ , RMSE and KGE) were calculated to compare the performance of the models during the testing (hold-out) period. The table support the selection of the best model for each study site between Random Forest (RF), eXtreme Gradient Boosting (XGBoost), Light Gradient Boosting (LGB), CatBoost (CTB), k-Nearest Neighbors (KNN), Support Vector Regression (SVR) and linear model. Overall, additionally, we tested other boosting methods, the best performance was metrics for testing phase can be found in the CatBoost model supplementary Table A2. Supplementary Table A1 contains the results during the training period for comparison.

Model/Metric	Lough Feeagh							Sau Reservoir						
	RF	XGB	LGB	CTB	LM	KNN	SVR	RF	XGB	LGB	CTB	LM	KNN	SVR
$R^2$	0.37	0.41	0.47	0.51	0.40	0.36	0.34	0.53	0.11	0.45	0.54	0.45	0.33	0.63
RMSE (QSU)	9.04	9.27	8.22	8.29	8.52	9.85	10.54	8.24	12.4	9.98	7.11	6.58	9.23	17.05
KGE	0.55	0.52	0.63	0.69	0.55	0.12	0.30	0.55	0.21	0.33	0.66	0.31	0.43	-0.82

### 249 3.3.1 Model prediction using all drivers compared to a reduced set

250 Figure 5 presents the fDOM prediction performance of the CatBoost model (best ML model overall) for Feeagh and Sau, using  
 251 different input configurations. The models were trained on 85% of the time series (blue points) and tested on the remaining 15%  
 252 using a hold-out approach (violet and green points). Two scenarios were compared: one using the most influential drivers (8  
 253 for Feeagh, 5 for Sau), and a second using only a reduced subset of reanalysis-based and easily accessible drivers, specifically  
 254 soil temperature, soil moisture, and Julian day.

255 For both lakes, the reduced driver models showed only a modest decline in predictive performance during the testing period.  
 256 For example, in Feeagh, the model using all influential drivers achieved  $R^2 = 0.51$  and  $KGE = 0.69$ , while the reduced model  
 257 still attained  $R^2 = 0.48$  and  $KGE = 0.67$ . Similarly, in Sau, the full model scored  $R^2 = 0.54$  and  $KGE = 0.66$ , whereas the reduced  
 258 model maintained a comparable  $R^2 = 0.50$  and  $KGE = 0.65$ . Although the training performance was higher in Sau compared  
 259 with the testing performance, indicating potential overfitting, the CatBoost model provided informative and generalisable  
 260 predictions in both study sites.

## 261 4 Discussion

### 262 4.1 Driver attribution in contrasting conditions

263 The most influential drivers for fDOM identified for each site suggest potential carbon-producing processes that drive DOM  
 264 dynamics in each lake. Interestingly, the most influential drivers were related to temperature and moisture at the deepest  
 265 soil layers, indicating a common relevance for carbon inputs from terrestrial processes in both catchments. The Irish site is  
 266 dominated by peat, a highly organic soil known to be sensitive to temperature-induced carbon release, especially as temperature  
 267 levels rise. This is supported by the relation of soil temperatures with fDOM at this site (Ryder et al., 2014). In contrast, water

268 availability constraints are much more pronounced in drier soils such as those of the Spanish site, as is validated by the partial  
269 dependence relationship of soil moisture with fDOM (Šimek et al., 2011). A plausible explanation is that, in Sau, increased soil  
270 moisture could boost biological activity, enhancing fDOM production, while in Feeagh, soil moisture showed a bell-shaped  
271 relationship, suggesting a more nuanced interplay between oxygen availability and microbial processes (Fig 4); [Fig. A9 and](#)  
272 [A10](#). In addition, higher soil moisture can reflect stronger soil–lake hydrological connectivity, particularly during wet periods  
273 and flushing events, facilitating the transport of terrestrially derived DOM into the lake.

274 DOM dynamics are driven by physical and biogeochemical processes in the soil that are sensitive to changes in temperature  
275 and moisture, e.g., microbial processes that break down organic matter (Kalbitz et al., 2000). The fact that the deepest layers  
276 of the soil were more important in the model for both sites than those shallower could be linked to potential carbon attenuation  
277 processes, such as soil organic matter decomposition and retention in the soil, including sorption processes (Dubeux et al.,  
278 2024; Rumpel and Kögel-Knabner, 2011). In any case, the production of carbon in the catchment that eventually ends up in  
279 the lake requires concurrent downstream transport, governed by rainfall events.

280 Climate and topography (see supplementary Fig. A1) dictate the flushing of accumulated DOM during rainfall events, but  
281 also can influence sustained baseflow DOM contributions. In Feeagh, carbon exports from the catchment have been observed  
282 regularly throughout the entire annual cycle, with a seasonal variability (Doyle et al., 2019). In contrast, DOM in Sau accumu-  
283 lates primarily during the summer and is mainly flushed out via surface runoff during the wetter winter months (Marcé et al.,  
284 2021). These patterns are supported by the relationship between inflow DOC concentration and fDOM at both sites (Fig. 4,  
285 [Figures A9 and A10](#)), which shows a slight increase in predicted fDOM under lower carbon input conditions. Thus, inflow  
286 DOC concentration could reflect discharge pulses and dilution effects driven by precipitation (Jennings et al., 2020), following  
287 the characteristic seasonality of each site.

288 Seasonality plays a crucial role in fDOM predictions, as evidenced by the relationship of the Julian day driver with DOM  
289 dynamics at both sites (Fig. 4). At the Irish site, DOM seasonality is primarily shaped by natural environmental processes,  
290 whereas in the Spanish site, human influence plays a much greater role. This distinction helps explain why surface lake water  
291 temperature and solar radiation, two variables typically linked to strong seasonal patterns, were important only for Feeagh,  
292 while reservoir volume was significant only for Sau ([Fig-Figures 3 and 4; Figures A9 and A10](#)). Volume and soil variables  
293 produce a similar effect on fDOM as Julian day at Sau, given that higher volumes closer to the winter season can lead to higher  
294 fDOM values. Incorporating Julian day into the workflow offers a simple yet effective way to represent seasonality, potentially  
295 replacing seasonal variables (e.g., air temperature) (see the correlation matrix of all drivers in supplementary Fig. [A5A12](#)).  
296 This proves that the use of machine learning approaches opens up opportunities to assess diverse drivers under contrasting  
297 conditions.

298 Improving the accuracy of DOM predictions in lakes can enhance efforts to reduce further water quality deterioration and  
299 support lake management. This study demonstrated a feasible approach for simulating daily fDOM in two contrasting lakes,  
300 especially when using the Catboost model given its good generalisability. The performance metrics (Fig. 5) obtained at each  
301 site for the different model simulations lie in a similar or better range than comparable studies that modelled carbon dynamics  
302 in lakes (e.g., Harkort and Duan, 2023; Liu et al., 2021; Zhang et al., 2021, 2024). It is important to be aware, however,

303 that previous studies are based on different frameworks. These include variations in the machine learning algorithms used,  
304 the target variable for quantifying carbon dynamics, with most studies having focused on DOC, whereas fDOM is the target  
305 variable here, as well as differences in input driver data and site-specific conditions. While both fDOM and DOC are widely  
306 used indicators of dissolved organic matter in lakes, they differ in measurement principles and in the fractions of organic  
307 matter they represent. Both variables, however, remain ecologically relevant for understanding DOM sources, transport, and  
308 transformation processes, which aligns with the context of this work.

## 309 **4.2 Scalability**

310 Our results suggest high potential for scalability, as predictive performance remained consistent across different driver sets  
311 even for two contrasting study sites, and performed good using only reanalysis data that is globally available and Julian day  
312 (Fig. 5). Importantly, this consistency remained even when specific highly-influential drivers were removed from the driver set.  
313 For instance, in Sau, where human intervention makes future reservoir outflows difficult to predict, avoiding reliance on water  
314 volume as a driver proved advantageous, as its removal from the set of drivers maintained model performance, despite being  
315 identified as an influential variable. It is likely that soil moisture at the deepest layer(see supplementary Fig. A5), a variable  
316 that showed a behavior similar to that of volume (Fig. 4), may have contributed to maintain the predictive performance when  
317 volume is removed from the driver set. In addition, for Feeagh, the predictive capacity was also maintained when using only  
318 meteorological and soil drivers. This demonstrates that a large driver dataset, such as the 24-variable set used in this study,  
319 would not be necessary to produce an accurate prediction when modelling fDOM using supervised machine learning.

## 320 **4.3 Limitations and future research**

321 Our approach offers the opportunity to validate and deploy a workflow capable of delivering daily DOM predictions in both  
322 undisturbed and anthropized sites, even when only limited data on input drivers are available, while at the same time providing  
323 insights into the dominant drivers. However, a site-specific model validation, including identification of appropriate training  
324 and testing periods, hyperparameter tuning for each specific study site and assessment of overfitting is essential. In terms of  
325 driver attribution, it is of note that the relationships identified using machine learning may not always be related at a process  
326 level (Sullivan, 2022). In our case, however, many of these same drivers had already been identified for river DOC levels in  
327 the Feeagh catchment (e.g., Doyle et al., 2019; Ryder et al., 2014). While the workflow can be easily replicated, fDOM data or  
328 data for another proxy for DOM are required. It is of note that such proxies of DOM are increasingly being incorporated into  
329 water quality monitoring programmes, an aspect that is convenient for testing workflows such as the one described (Downing  
330 et al., 2012).

331 The workflow presented here is not recommended for climate change studies, as the drivers of DOM variability can signifi-  
332 cantly change under entirely new and unrecorded climatic conditions. Consequently, supervised machine learning may fail to  
333 capture the signal from the time series. Moreover, the method's reliance on historical patterns limits its ability to extrapolate  
334 beyond the range of observed environmental conditions (Mi et al., 2024). Future research could expand the application of

335 this framework to a broader range of lakes, integrate additional drivers such as remote sensing-derived terrestrial and aquatic  
336 quantity and quality parameters (Duan et al., 2025).

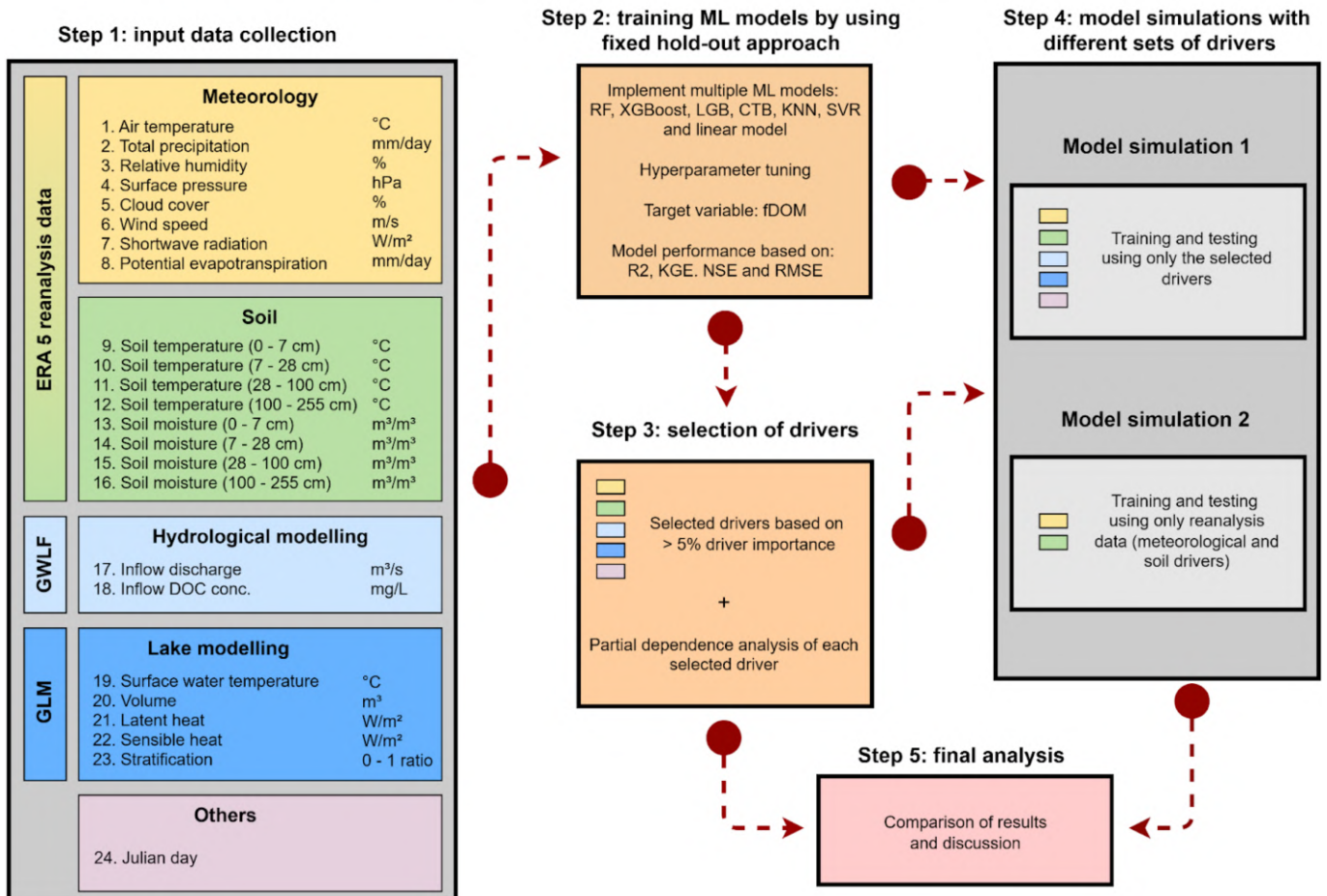
## 337 **5 Conclusions**

338 By identifying the key environmental drivers of lake dissolved organic matter (DOM) dynamics, this study presents an open,  
339 robust and scalable workflow for daily DOM prediction using different ML algorithms. Validated in two hydroclimatic con-  
340 trasting sites in Ireland (Lough Feeagh) and Spain (Sau Reservoir), the approach revealed that deep soil temperature is the  
341 dominant driver in the peat-rich, temperate Irish catchment, whereas deep soil moisture plays a more critical role in the drier,  
342 Mediterranean setting of the Spanish site. These primary drivers are further shaped by hydrological processes, seasonal vari-  
343 ability, and human activities.

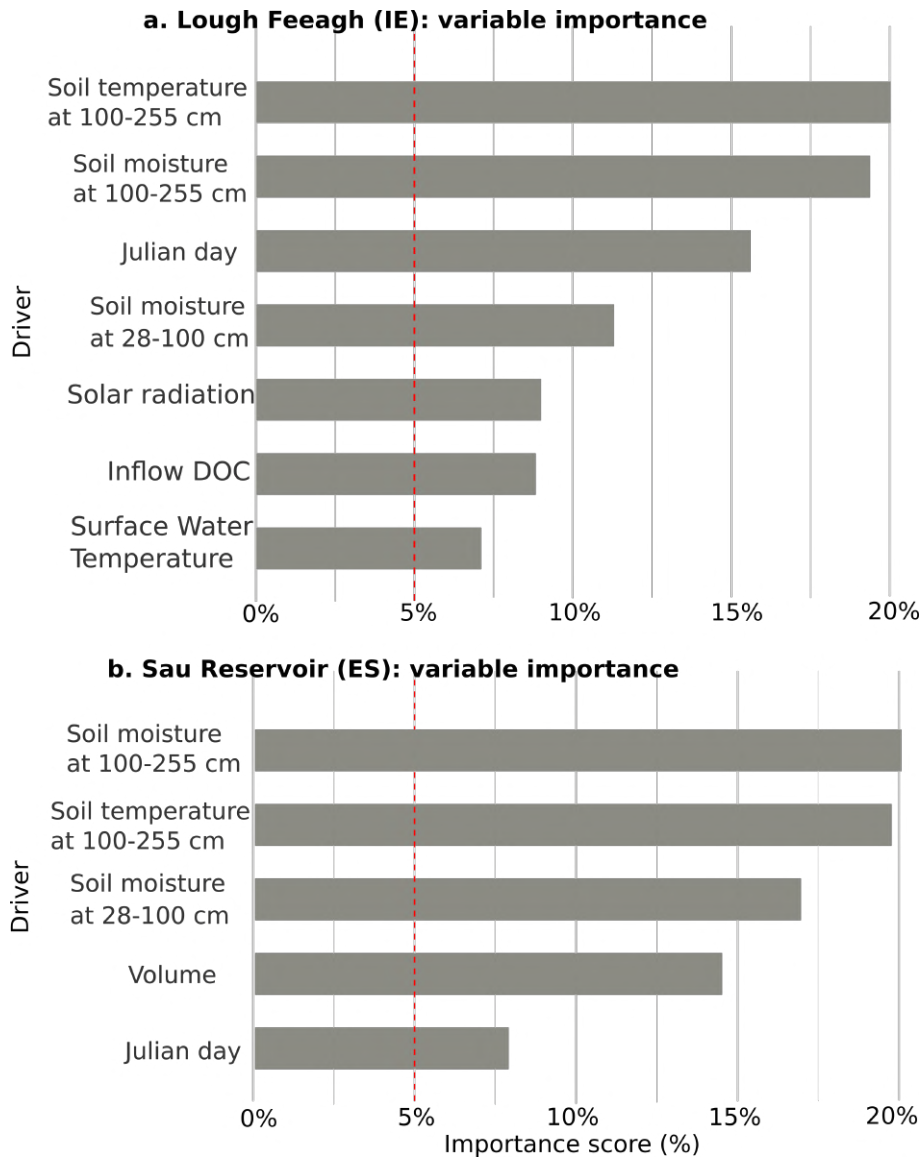
344 The workflow showed good predictive performance even when based solely on globally available reanalysis data, supporting  
345 its potential applicability to other freshwater systems worldwide. In addition to expanding the set of approaches available for  
346 lake DOM prediction, the workflow offers transparent driver attribution, contributing valuable insights into the natural and  
347 anthropogenic processes governing carbon cycling in aquatic ecosystems.

348 *Code and data availability.* All data and codes used in this study are available in this repository: [https://github.com/danielmerbet/driver\\_](https://github.com/danielmerbet/driver_)  
349 [attribution\\_fdom](https://github.com/danielmerbet/driver_attribution_fdom). A DOI will be provided in the final version of the manuscript

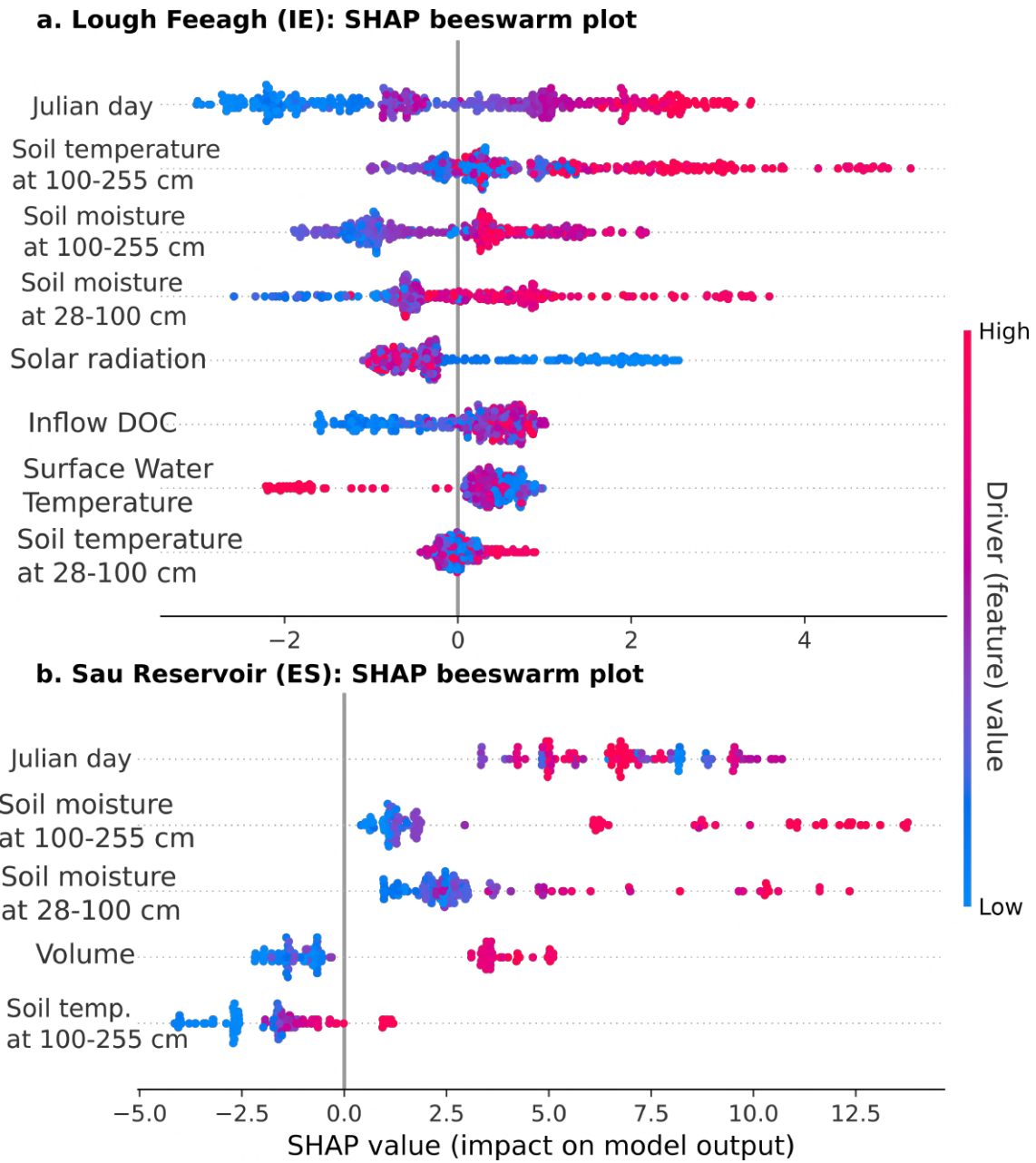
## Workflow



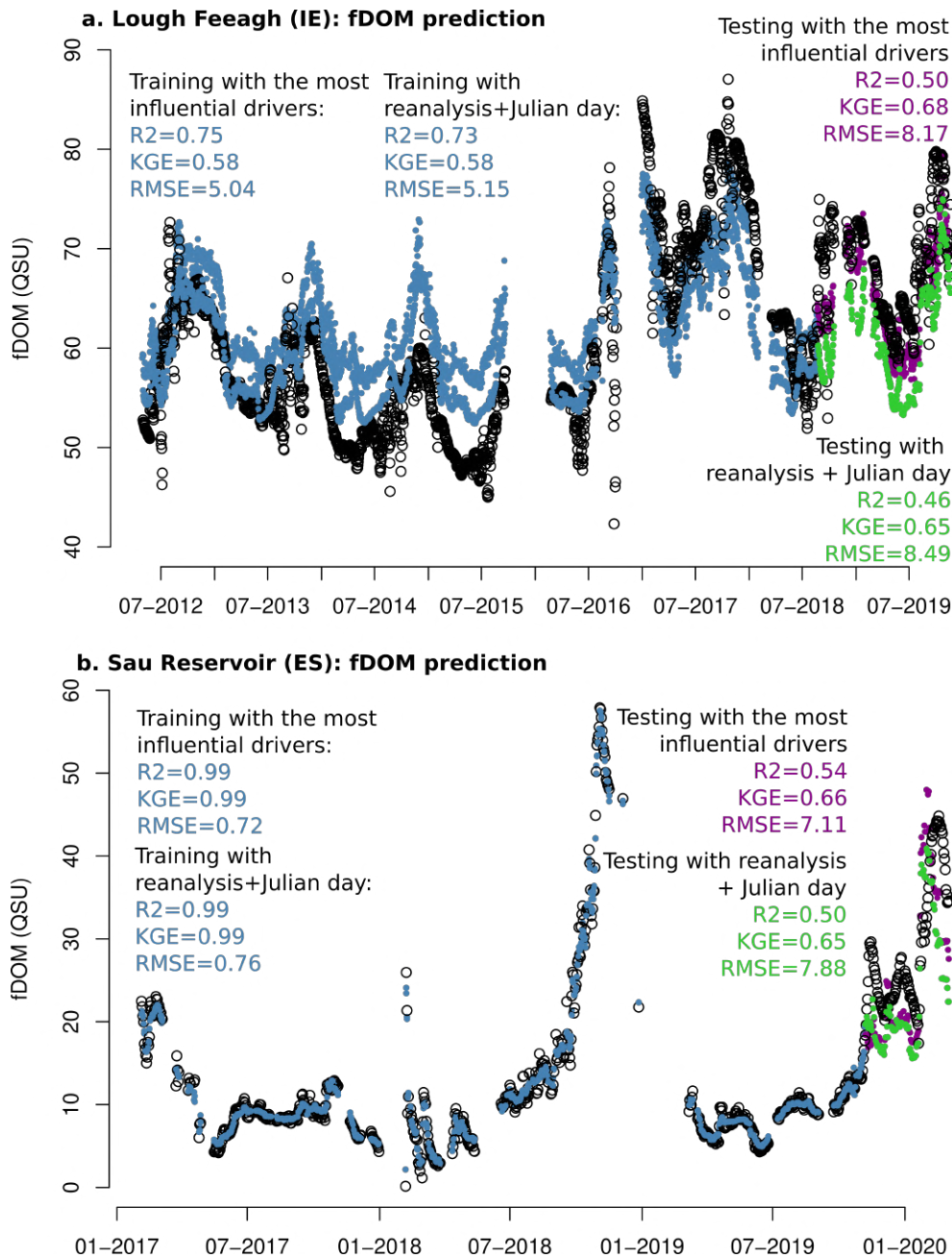
**Figure 2.** Workflow implemented for obtaining fDOM predictions in both study sites. The process consists of workflow includes five steps: (1) Collecting input data representing all potential drivers compilation of fDOM, input variables including climate data meteorology (yellow) and soil data (green) from globally accessible reanalysis data from ERA5 (the fifth-generation atmospheric reanalysis product produced by ECMWF), hydrologic model The Generalised Watershed Loading Functions Model (GWLF) outputs (light blue), lake model The General Lake Model (GLM) outputs (blue), and other sources Julian day (light magenta); (2) Training the ML models by splitting the data inspection and model training using a split between training (85% of the time series) and testing (15% of the time series) periods; (3) Selecting key drivers by assessing their contribution to node purity or gain contribution in the ML models, only drivers exceeding 5% selection of variable importance were retained for partial dependence analyses; influential drivers based on feature importance metrics; (4) Running two model simulations using only the most influential selected drivers, and from these using only a reduced subset based on globally accessible available reanalysis drivers variables and Julian day for ease; and scalable implementation; (5) Analyzing and comparing the modelling results comparison of model outputs. DOC stands for Dissolved Organic Carbon. The workflow is available at: [https://github.com/danielmerbet/driver\\_attribution\\_fdom](https://github.com/danielmerbet/driver_attribution_fdom)



**Figure 3.** ~~Selecting Selection of~~ influential drivers ~~to predict for predicting~~ fDOM. ~~24 Twenty-four~~ drivers from ~~various multiple~~ sources were used to train the ~~four ML machine-learning~~ models that directly provide feature importance: Random Forest (RF), eXtreme Gradient Boosting (XGBoost), Light Gradient Boosting (LGB), CatBoost (CTB). ~~These included 8 including~~ climate variables, ~~8 soil variables,~~ ~~2 outputs from hydrologic and water quality modelling,~~ ~~5 water quality model outputs from,~~ lake modelling model outputs, and the cosine of the cosine-transformed Julian day (see Figure 2 for more details). ~~The most relevant drivers were identified based on their contribution to node purity or Feature importance is shown for the best-performing model (CatBoost) using~~ gain contribution, ~~with only those and drivers~~ exceeding a 5% variable importance being selected. For both case studies, the key drivers were soil temperature at 100–255 cm, soil moisture at 28–100 cm and 100–255 cm, and Julian day. Additionally, lake surface water temperature, inflow DOC concentration, and soil temperature at 28–100 cm were selected for Feeagh, while lake volume was selected for Saure highlighted.



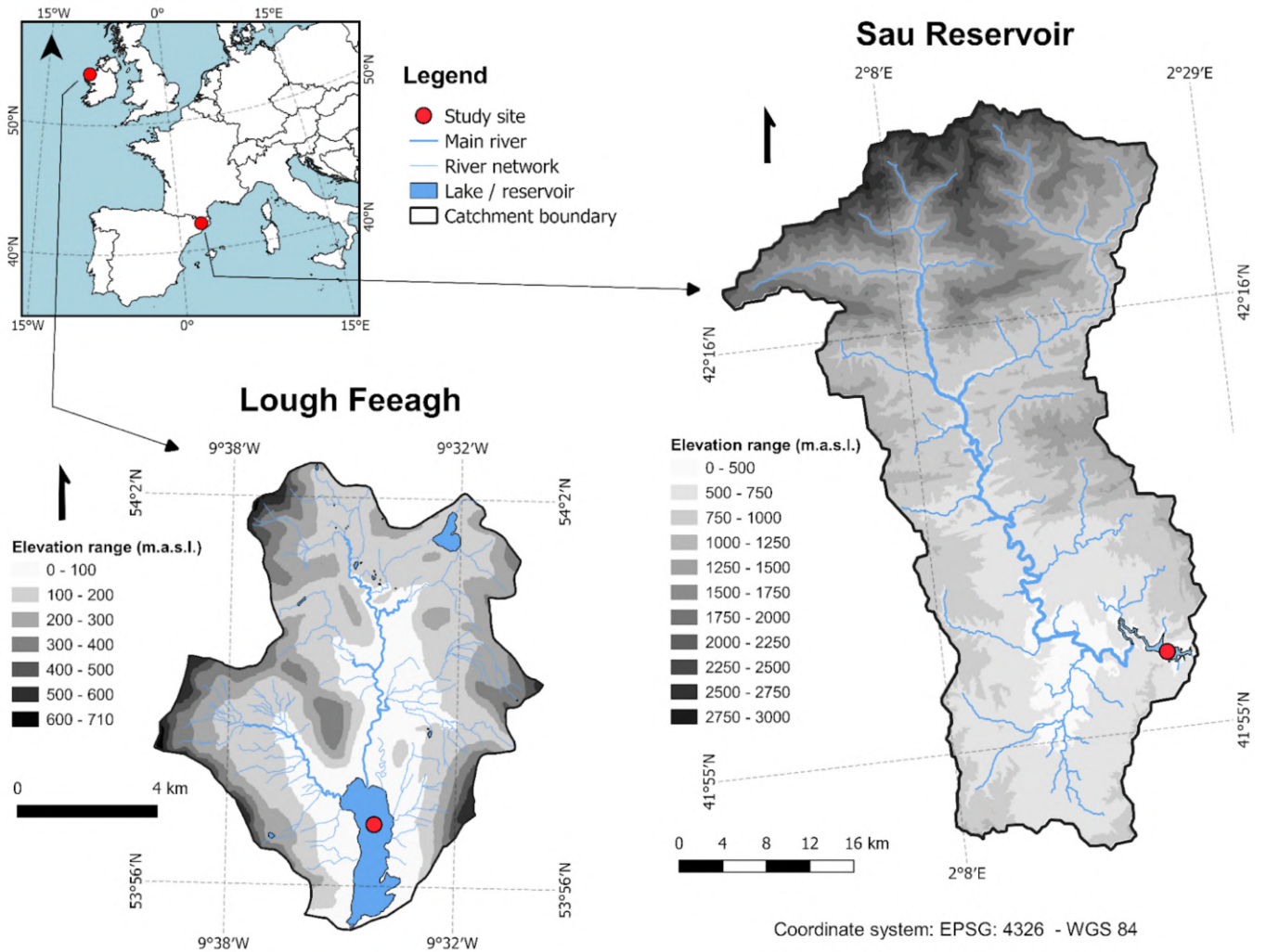
**Figure 4.** Driver influence in predicting fDOM. It presents the partial dependence plots for (a.) Lough Feeagh and (b.) Sau Reservoir based on the random forest (RF) model, and SHAP beeswarm plots showing the contribution of key drivers to fDOM predictions for (e.a) Lough Feeagh and (d.b) Sau Reservoir based on the CatBoost model (CTB) model. In Feeagh, soil temperature at the deepest layer strongly influenced fDOM predictions, while in Sau, soil moisture at the deepest layer plays a significant role due to water availability constraints. In both cases, increases in these key drivers correspond to increases in fDOM. In both study sites Points represent individual observations, Julian day was a relevant factor coloured by feature value, although the influence of seasonality on fDOM predictions was more evident in Feeagh and are ordered by mean absolute SHAP value.



**Figure 5.** fDOM predictions obtained from the CatBoost model using different driver sets. Here the Panels show training (blue, 85% of the time series) and testing periods (violet or green, 15% of the time series) periods are shown. The training and testing periods are completely independent and follow the hold-out period method, and model performance metrics ( $R^2$ , KGE, and RMSE) are calculated for both periods:

a. Training and testing results for (a) Feagh: simulations using the 8-most influential drivers (violet) and using only easily-accessible and reanalysis data (soil temperature, soil moisture and a reduced subset of reanalysis-based drivers + Julian day) from those influential drivers (green);

b. Training and testing results for (b) Sau: simulations using the 5-most influential drivers (violet) and using only easily-accessible and reanalysis data (soil temperature, soil moisture and a reduced subset of reanalysis-based drivers + Julian day) from those influential drivers (green). Model performance metrics ( $R^2$ , KGE, RMSE) are reported for both periods.



**Figure A1.** Elevation range for the two contrasting freshwater ecosystems: Lough Feeagh and Sau Reservoir.

350 **Appendix A: Supplementary information**

351 **Topography of Feeagh and Sau catchments**

352 Figure A1 contains the elevation range for the two contrasting freshwater ecosystems: Lough Feeagh and Sau Reservoir.

353 ~~fDOM-correction~~ CORINE land-cover classes and aggregation for both study sites catchments

354 Here, we present Figure A2 listing all original CORINE Land Cover 2018 classes identified within each catchment and their  
355 aggregation into the grouped categories used in Figure 1. These tables show how individual CORINE classes were consolidated  
356 to represent dominant catchment characteristics while maintaining consistency with the original CORINE classification. This  
357 information supports the visual interpretation presented in Figure 1.

a.

Lough Feeagh			
Code	Colour	CODE_18	CORINE Land Cover Classification
312		312	Coniferous forest
321		321	Natural grassland
322		322	Moors and heathland
324		324	Transitional woodland/shrubs
333		333	Sparsely vegetated areas
412		412	Peatbogs
511		511	Water courses
512		512	Water bodies

Land cover group	% of total area
Organic-rich soil LC	45%
Agriculture	0%
Grassland	28%
Urban	0%
Forest	25%
Watercourse	2%
	100%

b.

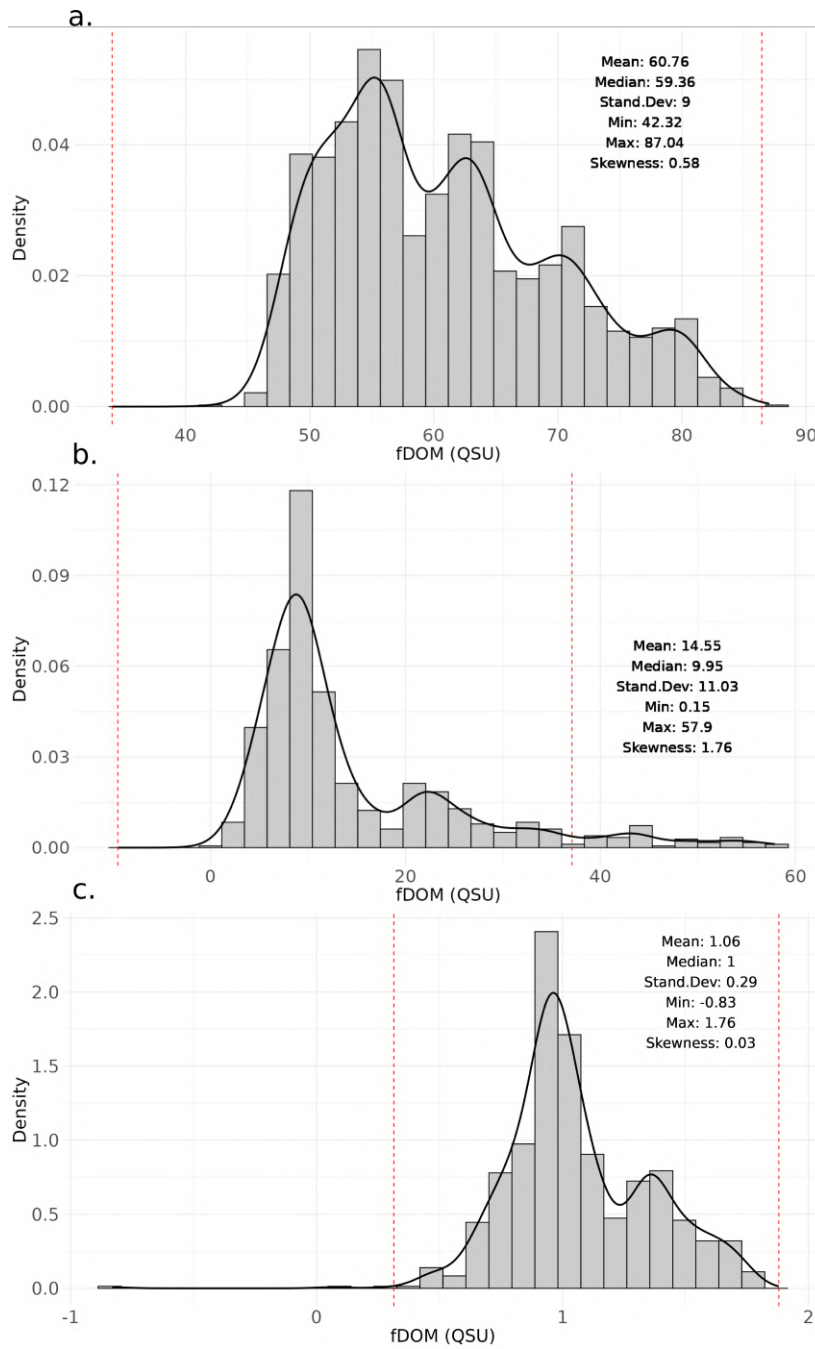
Sau Reservoir			
Code	Colour	CODE_18	CORINE Land Cover Classification
111		111	Continuous urban fabric
112		112	Discontinuous urban fabric
121		121	Industrial or commercial units, public services and military installations
131		131	Mineral extraction sites
142		142	Green urban areas
211		211	Non-irrigated arable land
212		212	Permanently irrigated arable land *
231		231	Pastures
242		242	Complex cultivation patterns
243		243	Land principally occupied by agriculture
311		311	Broad-leaved forest
312		312	Coniferous forest
313		313	Mixed forest
321		321	Natural grassland
322		322	Moors and heathland
323		323	Sclerophyllous vegetation*
324		324	Transitional woodland/shrubs
332		332	Bare rock
333		333	Sparsely vegetated areas
511		511	Water courses
512		512	Water bodies

Land cover group	% of total area
Organic-rich soil LC	6%
Agriculture	18%
Grassland	12%
Urban	4%
Forest	59%
Watercourse	0%
	100%

Figure A2. a. CORINE Land Cover 2018 classes identified in Feeagh and the grouped categories shown for visual clarity in the main manuscript in Figure 1; b. CORINE Land Cover 2018 classes identified in Sau and the grouped categories shown for visual clarity in the main manuscript in Figure 1.

358 **Data inspection and exploration**

359 Data inspection and exploration: the target variable (fDOM) at one of the study sites (Feeagh) exhibits low to moderate  
360 skewness (approximately 0.5) with relatively few outliers, whereas the other study site (Sau) displays higher skewness (greater  
361 than 1 but below 2) and several extreme values. Zero-inflation does not appear to be an issue at either site.



**Figure A3.** Exploratory data analysis for fDOM in both study sites, a. Shows the density plot of fDOM in Feeagh and b. Shows the density plot of fDOM in Sau, c. Shows the density plot of fDOM in Sau using log10 transformation. The red vertical line separates the outliers.

362 Water temperature correction of fDOM for both study sites

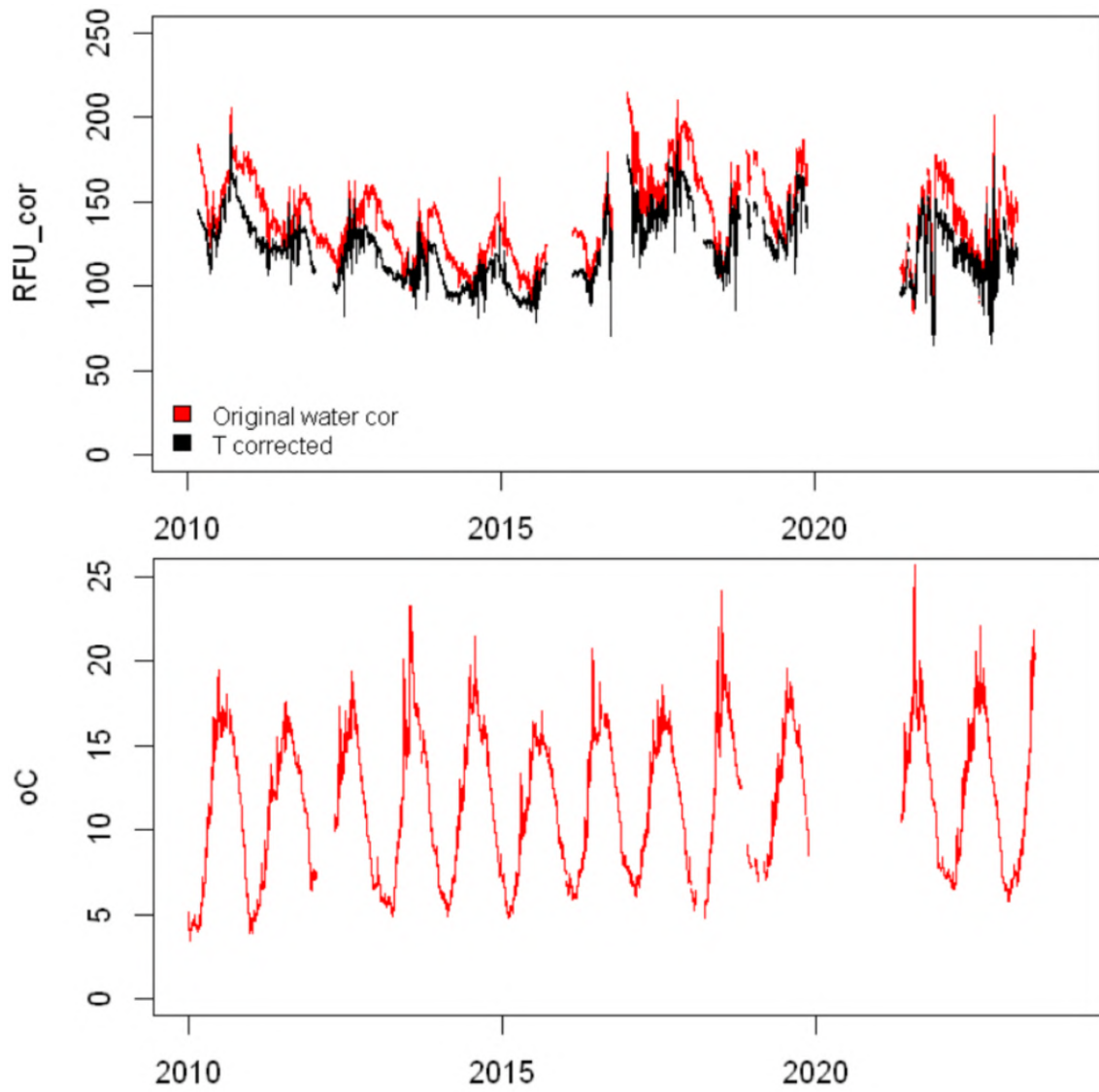
In Feeagh, fDOM data was corrected for the temperature quenching effect in previous scientific studies (Doyle et al., 2019; Ryder et al., 2011).

The raw measurements were corrected using the compensation approach described by (Watras et al., 2011; Ryder et al., 2012)

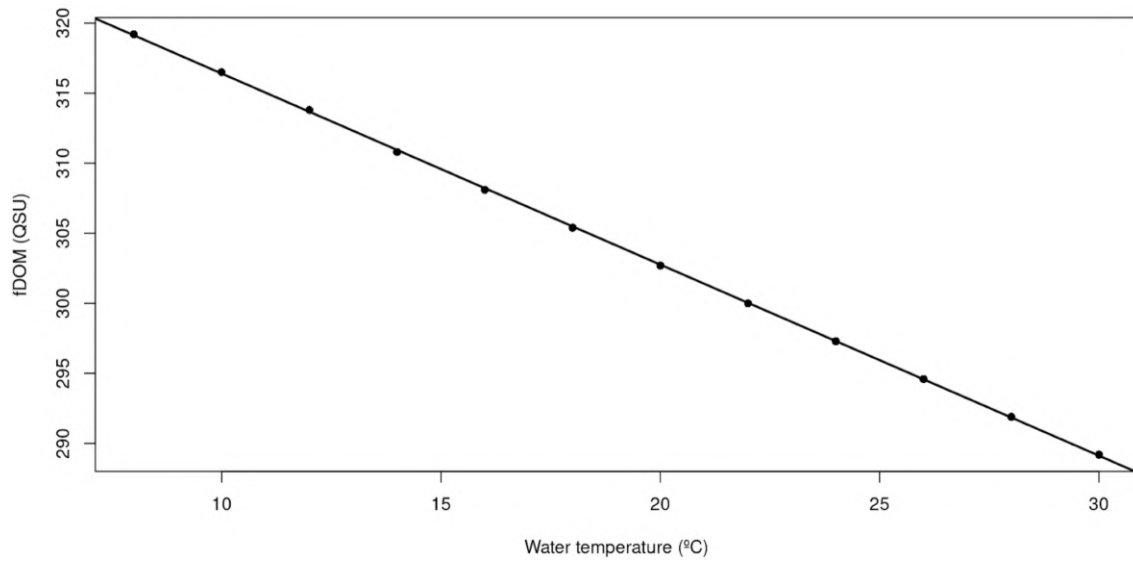
Corrected fluorescence was calculated relative to a reference temperature of 20°C using: Corrected fDOM =  $\frac{\text{raw fDOM}}{1+(T-20)^k}$  1000

363 Where T is the sonde water temperature (°C) and k is the temperature correction coefficient (k = 0.0168 for Feeagh, derived

364 as the average of monthly values reported by Ryder et al. (2012)).

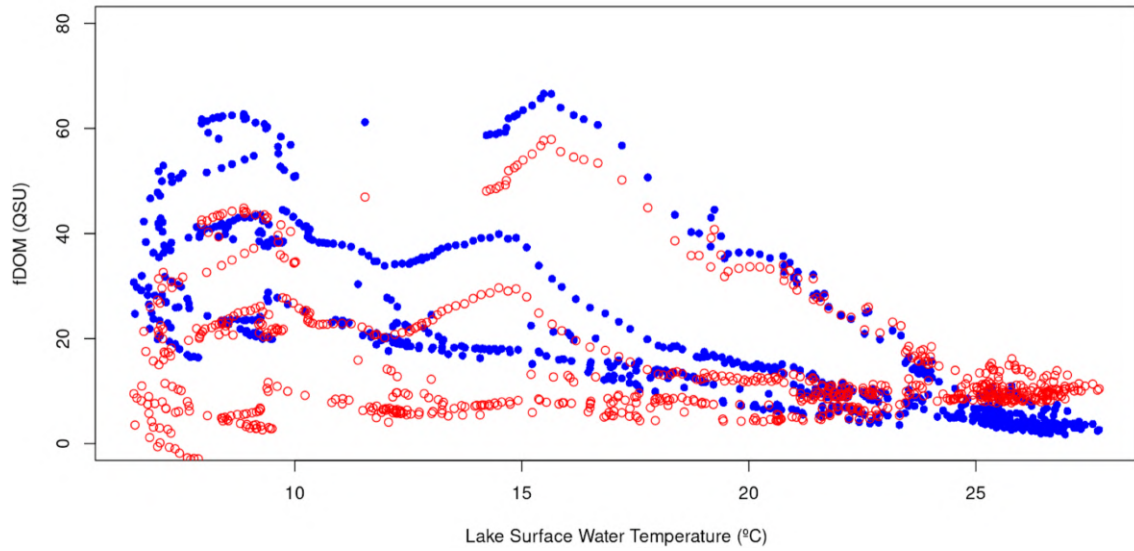


**Figure A4.** Upper panel: Raw (blue) and temperature-corrected (red) fDOM measurements for Feeagh; lowe panel: surface water temperature sonde measurements at Feeagh.



**Figure A5.** fDOM relation with temperature with a sample of 300 QSU provided by the manufacturer of the fDOM sensor in Sau.

365 In Sau, the fDOM data were corrected for the temperature quenching effect, following a test provided by the fDOM sensor  
366 manufacturer, where they use a 300 QSU sample of water and change the temperature to get the effect of temperature in the  
367 measurement, results can be found in [Figure A2A5](#).



**Figure A6.** Uncorrected (blue) and corrected (red) fDOM data for Sau.

368 Then, fDOM data in Sau was corrected following this linear regression and surface water temperature on the lake. Figure [A3](#)  
369 [A6](#) presents the uncorrected values in blue and corrected values in red (8 negative values were removed from the total sample  
370 of 777)

## 371 **Hydrologic Modelling**

372 Daily time series of inflow discharge and inflow DOC concentration into each site were generated using the Generalised Wa-  
373 tershed Loading Functions Model (GWLF) coupled with a DOC module (GWLF-DOC). This model simulates catchment  
374 hydrology (water balance and water distribution among the different hydrological pathways) and DOC dynamics (DOC pro-  
375 duction and DOC washout) in a daily time step. The model input requirements include daily time series of two meteorological  
376 variables: total precipitation and air temperature; as well as land cover, land use, and soil characterisation.

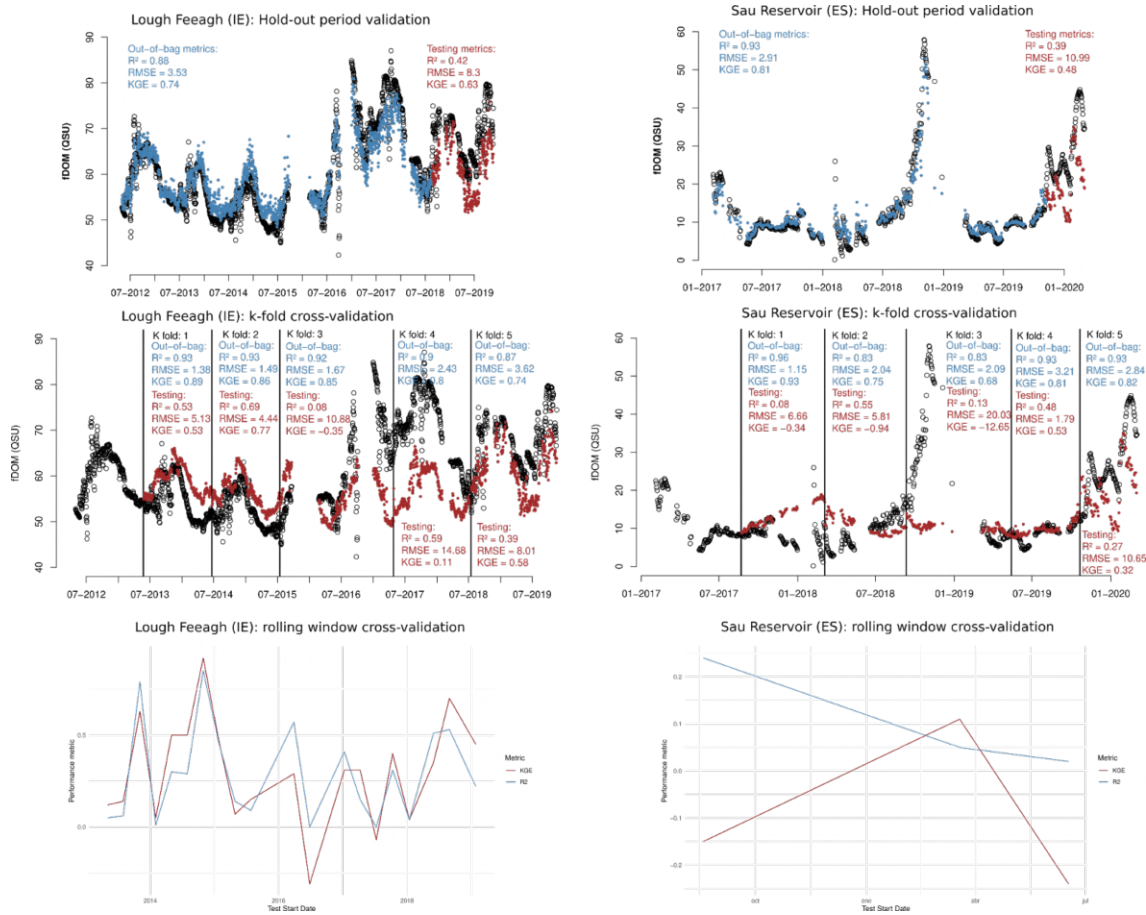
377 GWLF-DOC was applied to Feeagh based on previous model applications in the Irish catchment (Paíz et al., 2025a), for  
378 which measured discharge data were used to calibrate and validate the hydrology (2013-2018 and 2019-2023, respectively),  
379 and DOC concentration data were used to calibrate the DOC module (2016-2023). In Sau, observed inflow discharge data  
380 were used to calibrate and validate the hydrology (2008-2011 and 2011-2024, respectively), and measured DOC concentration  
381 data were used to calibrate and validate the DOC module (2008-2014 and 2016-2018, respectively) using the same calibration  
382 strategy than for the Irish site. Calibration results were satisfactory for both hydrology (Feeagh:  $R^2 = 0.64$  and  $NSE = 0.64$ ;  
383 Sau:  $R^2 = 0.66$  and  $NSE = 0.66$ ;) and DOC (Feeagh:  $R^2 = 0.45$  and  $NSE = 0.47$ ; Sau:  $R^2 = 0.44$  and  $NSE = 0.40$ ). Similarly,  
384 validation results were satisfactory for both hydrology (Feeagh:  $R^2 = 0.60$  and  $NSE = 0.60$ ; Sau:  $R^2 = 0.42$  and  $NSE = 0.42$ )  
385 and DOC in the case of Sau ( $R^2 = 0.50$  and  $NSE = 0.46$ ).

## 386 **Lake Modelling**

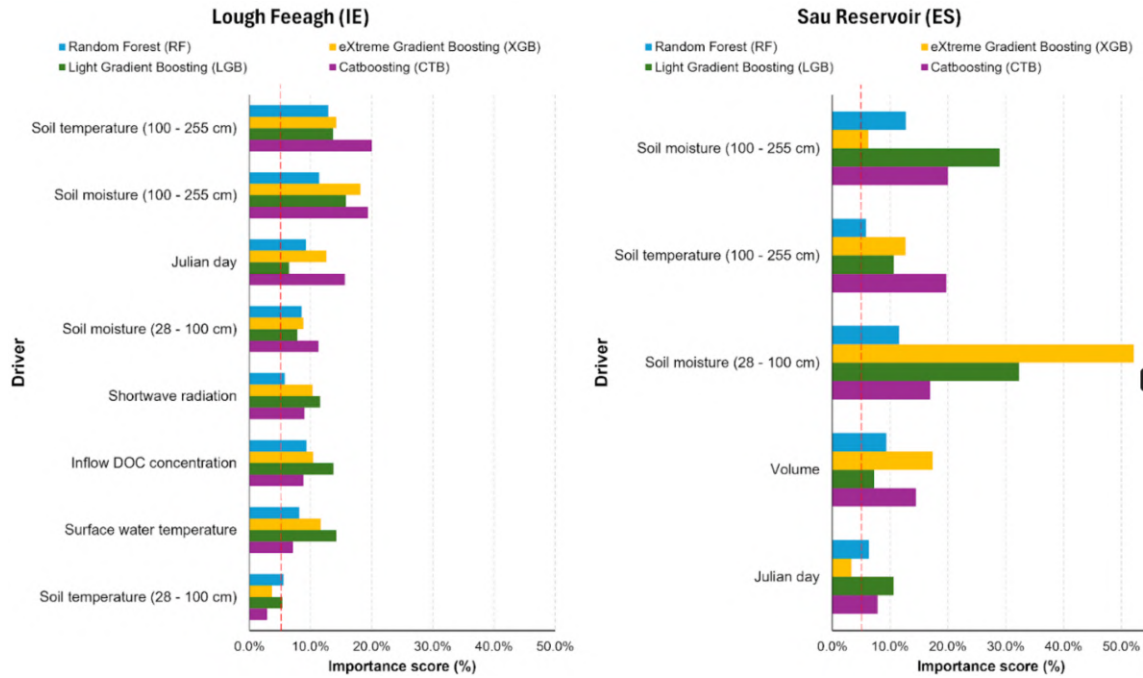
387 Daily time series of 5 key lake variables (see Fig. 2) were obtained from the General Lake Model (GLM) run for each site. GLM  
388 is an open-source, one-dimensional hydrodynamic model designed to simulate the vertical stratification and water balance of  
389 lakes and reservoirs. It calculates vertical profiles of temperature, and density by accounting for factors such as inflows and out-  
390 flows, mixing processes, and surface heating and cooling (Hipsey et al., 2019). GLM was calibrated and validated by evaluating  
391 the fit of modelled water temperature against measured water temperature profile data in Feeagh (2010-2015 and 2016-2017,  
392 respectively) and Sau (1997-2007 and 2008-2018, respectively). The calibration strategy was based on previous lake modelling  
393 deployments at each site (~~Mercado-Bettín et al., 2021; Paíz et al., 2025a~~)([Mercado-Bettín et al., 2021; Paíz et al., 2025b](#)). Model  
394 performance was satisfactory for both sites.

## 395 **Comparison of validation methods using Random Forest**

396 To pick the most suitable validation method, we implemented hold-out period method used in the main manuscript, k-fold  
397 cross-validation method using  $k=5$ , and rolling window cross-validation using training size of two year, testing size of one  
398 year, and a shift window every 90 days. Results are shown in Figure [A4A7](#).

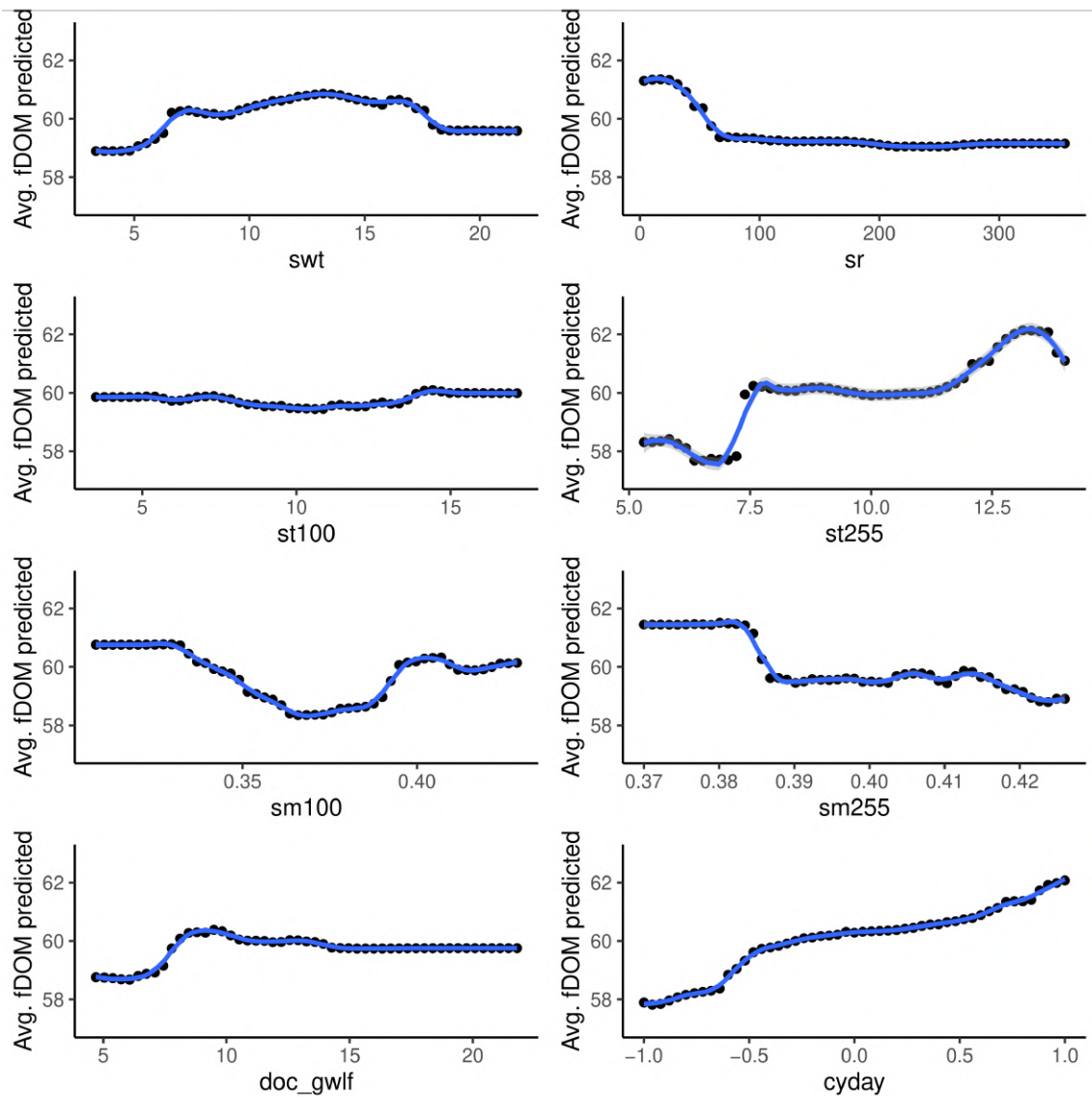


**Figure A7.** Comparison of validation methods using random forest for both study sites: hold-out period method used in the main manuscript, k-fold cross-validation method using  $k=5$ , and rolling window cross-validation using training size of two year, testing size of one year, and a shift window every 90 days. For this method, in the case of Sau Reservoir it is not possible to get a clear analysis due to the limited data.



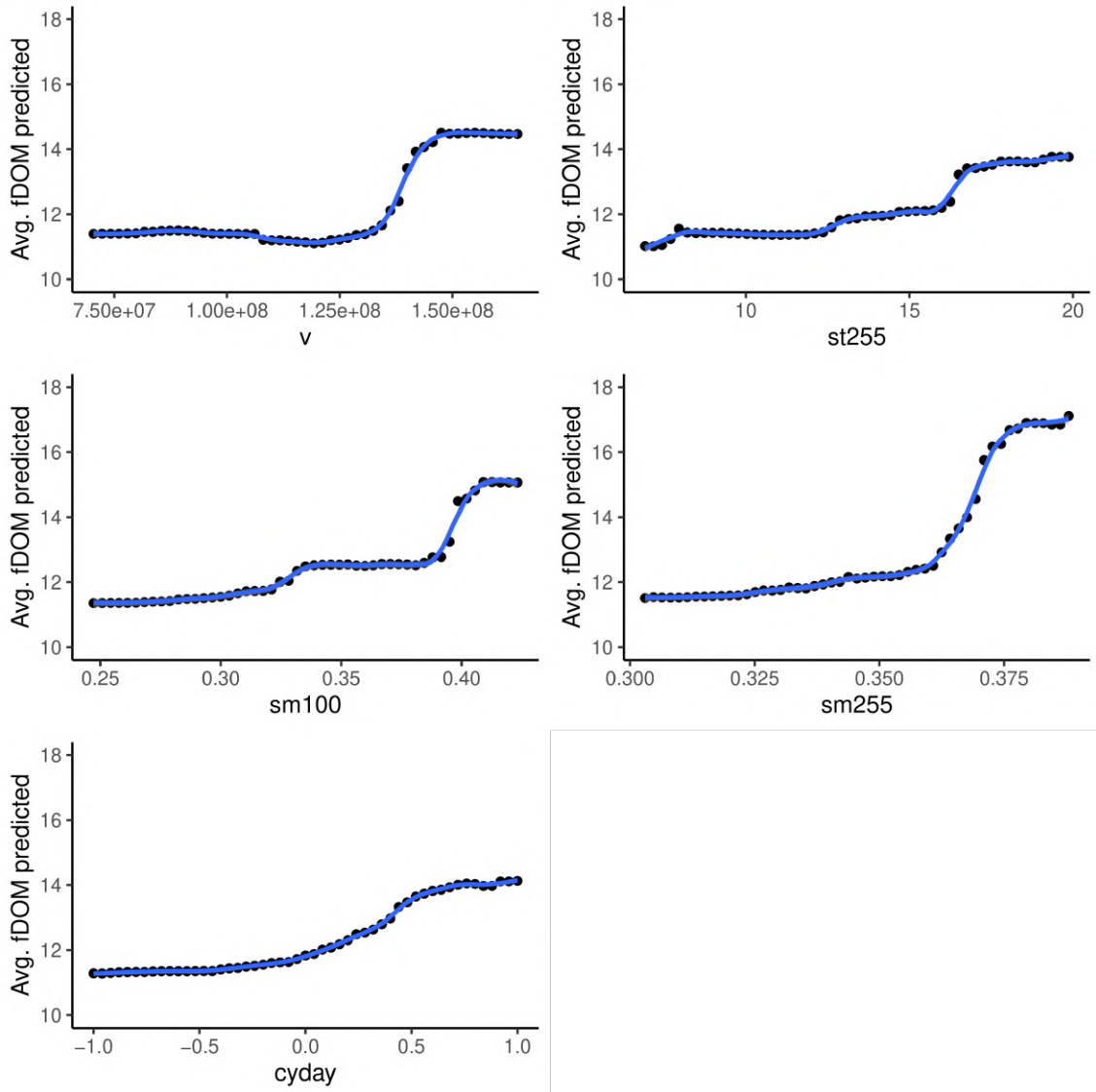
**Figure A8.** Selection of influential drivers for predicting fDOM. Twenty-four drivers from multiple sources were used to train the machine-learning models, including climate variables, soil variables, hydrologic and water-quality model outputs, lake model outputs, and cosine-transformed Julian day. Feature importance is shown for four models Random Forest, Light Gradient Boosting, eXtreme Gradient Boosting and Catboosting, and drivers exceeding 5% importance are highlighted.

a. Lough Feeagh (IE)

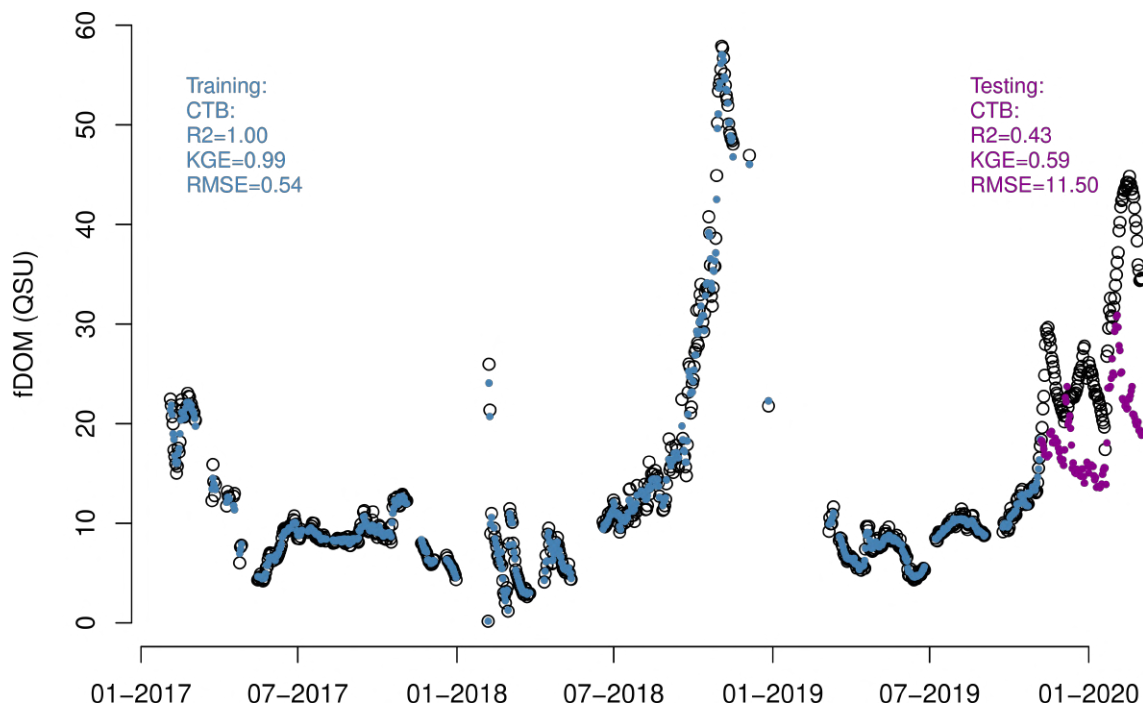


**Figure A9.** Partial dependence plots for the most influential drivers (feature importance > 5%, estimated using node purity) based on Random Forest model outputs for Lough Feeagh. Cosine-transformed Julian day (cyday) values are shown for seasonality; approximate calendar timing corresponds to cos(Julian day) 1 (winter), 0 (spring/autumn), and 1 (summer).

b. Sau Reservoir (ES)



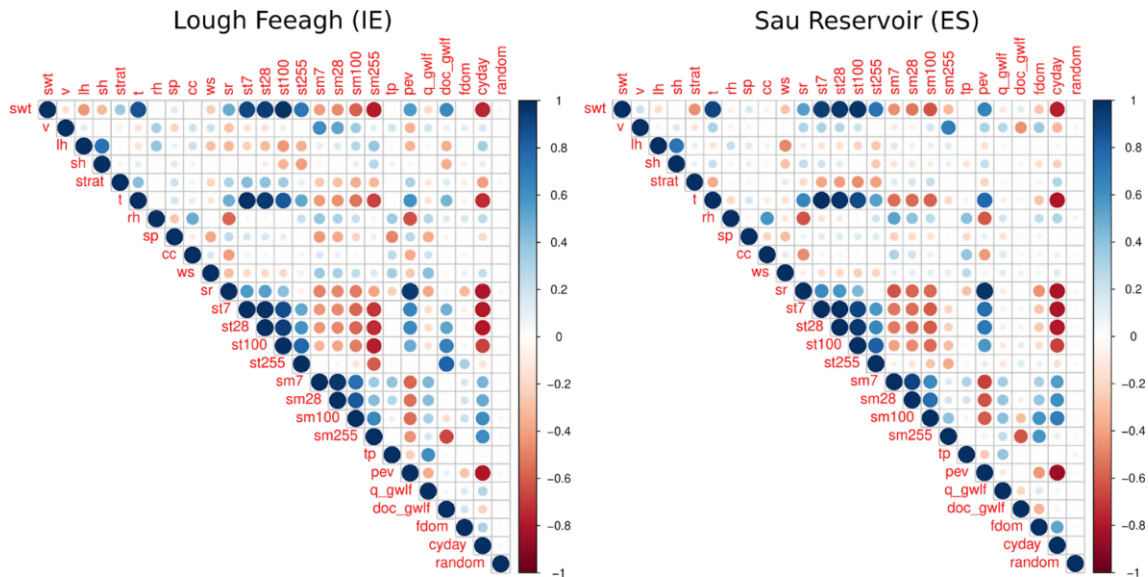
**Figure A10.** Partial dependence plots for the most influential drivers (feature importance > 5%, estimated using node purity) based on Random Forest model outputs for Sau Reservoir. Cosine-transformed Julian day (cyday) values are shown for seasonality; approximate calendar timing corresponds to  $\cos(\text{Julian day}) = 1$  (winter), 0 (spring/autumn), and 1 (summer).



**Figure A11.** fDOM predictions obtained from the boosting (CatBoost) model using the most influential drivers with initial log10 transformation of the fDOM data.

402 **Correlation matrix of all drivers and fDOM for Feeagh and Sau.**

Correlation matrix of all drivers is presented in Figure [A5A12](#).



**Figure A12.** Correlation matrix of all drivers and fDOM for Feeagh and Sau.

403

404 **Model comparison during training** to predict fDOM using **all selected the most influential drivers** in both study sites.

405 Resulting metrics during **training period testing and training periods** for all models are shown in Table A1 [and A2, respectively](#).

**Table A1.** [Model comparison during testing to predict fDOM using all selected drivers in both study sites. Statistic metrics \( \$R^2\$ , RMSE and KGE\) were calculated to compare the performance of the models Random Forest \(RF\), eXtreme Gradient Boosting \(XGBoost\), Light Gradient Boosting \(LGB\), Catboosting \(CTB\), k-Nearest Neighbors \(KNN\), Support Vector Regression \(SVR\) and linear model during training period.](#)

	Lough Feeagh							Sau Reservoir						
<u>Model/Metric</u>	<u>RF</u>	<u>XGB</u>	<u>LGB</u>	<u>CTB</u>	<u>LM</u>	<u>KNN</u>	<u>SVR</u>	<u>RF</u>	<u>XGB</u>	<u>LGB</u>	<u>CTB</u>	<u>LM</u>	<u>KNN</u>	<u>SVR</u>
<u><math>R^2</math></u>	<u>0.37</u>	<u>0.40</u>	<u>0.40</u>	<u>0.50</u>	<u>0.40</u>	<u>0.36</u>	<u>0.34</u>	<u>0.53</u>	<u>0.09</u>	<u>0.25</u>	<u>0.54</u>	<u>0.45</u>	<u>0.33</u>	<u>0.63</u>
<u>RMSE</u>	<u>9.07</u>	<u>9.28</u>	<u>8.97</u>	<u>8.17</u>	<u>8.52</u>	<u>9.85</u>	<u>10.54</u>	<u>8.51</u>	<u>12.22</u>	<u>11.35</u>	<u>7.11</u>	<u>6.58</u>	<u>9.23</u>	<u>17.05</u>
<u>KGE</u>	<u>0.54</u>	<u>0.51</u>	<u>0.50</u>	<u>0.68</u>	<u>0.55</u>	<u>0.12</u>	<u>0.30</u>	<u>0.55</u>	<u>0.25</u>	<u>0.22</u>	<u>0.66</u>	<u>0.31</u>	<u>0.43</u>	<u>-0.82</u>

**Table A2.** Model comparison during training to predict fDOM using all selected drivers in both study sites. Statistic metrics ( $R^2$ , RMSE and KGE) were calculated to compare the performance of the models Random Forest (RF), eXtreme Gradient Boosting (XGBoost), Light Gradient Boosting (LGB), Catboosting (CTB), k-Nearest Neighbors (KNN), Support Vector Regression (SVR) and linear model during training period.

Model/Metric	Lough Feeagh							Sau Reservoir						
	RF	XGB	LGB	CTB	LM	KNN	SVR	RF	XGB	LGB	CTB	LM	KNN	SVR
$R^2$	0.99	1.00	<del>0.88</del> <u>0.99</u>	<del>0.71</del> <u>0.75</u>	0.22	0.99	0.89	0.99	1.00	<del>0.98</del> <u>1.00</u>	0.99	0.67	0.98	0.89
RMSE (QSU)	0.97	<del>0.40</del> <u>0.00</u>	<del>3.35</del> <u>0.86</u>	<del>5.19</del> <u>5.04</u>	7.87	0.92	2.95	0.81	<del>0.11</del> <u>0.01</u>	<del>1.48</del> <u>0.17</u>	0.72	5.43	1.17	1.17
KGE	0.95	<del>0.99</del> <u>1.00</u>	<del>0.79</del> <u>0.98</u>	<del>0.59</del> <u>0.58</u>	0.25	0.98	0.90	0.97	1.00	<del>0.97</del> <u>1.00</u>	0.99	0.75	0.99	0.90

406 Model comparison to predict fDOM using the most influential drivers in Sau with previous application of log10  
407 transformation.

408 Tables A3 and A4 presents the testing and training performance metrics results for Sau site with previous application of log10  
409 transformation.

**Table A3.** Model comparison during testing to predict fDOM using the most influential drivers in Sau with previous application of log10 transformation. Statistic metrics ( $R^2$ , RMSE and KGE) were calculated to compare the performance of the models Random Forest (RF), eXtreme Gradient Boosting (XGBoost), Light Gradient Boosting (LGB), Catboosting (CTB), k-Nearest Neighbors (KNN), Support Vector Regression (SVR) and linear model during training period.

Sau Reservoir							
<u>Model/Metric</u>	<u>RF</u>	<u>XGB</u>	<u>LGB</u>	<u>CTB</u>	<u>LM</u>	<u>KNN</u>	<u>SVR</u>
<u><math>R^2</math></u>	<u>0.51</u>	<u>0.18</u>	<u>0.22</u>	<u>0.43</u>	<u>0.40</u>	<u>0.26</u>	<u>0.58</u>
<u>RMSE (QSU)</u>	<u>14.24</u>	<u>16.56</u>	<u>13.07</u>	<u>11.50</u>	<u>10.91</u>	<u>9.75</u>	<u>35.23</u>
<u>KGE</u>	<u>0.63</u>	<u>0.38</u>	<u>0.34</u>	<u>0.59</u>	<u>0.36</u>	<u>0.32</u>	<u>-1.35</u>

**Table A4.** Model comparison during training to predict fDOM using the most influential drivers in Sau with previous application of log10 transformation. Statistic metrics ( $R^2$ , RMSE and KGE) were calculated to compare the performance of the models Random Forest (RF), eXtreme Gradient Boosting (XGBoost), Light Gradient Boosting (LGB), Catboosting (CTB), k-Nearest Neighbors (KNN), Support Vector Regression (SVR) and linear model during training period.

Sau Reservoir							
<u>Model/Metric</u>	<u>RF</u>	<u>XGB</u>	<u>LGB</u>	<u>CTB</u>	<u>LM</u>	<u>KNN</u>	<u>SVR</u>
<u><math>R^2</math></u>	<u>0.98</u>	<u>1.00</u>	<u>0.99</u>	<u>1.00</u>	<u>0.58</u>	<u>0.94</u>	<u>0.87</u>
<u>RMSE (QSU)</u>	<u>1.05</u>	<u>0.07</u>	<u>0.74</u>	<u>0.54</u>	<u>4.65</u>	<u>1.22</u>	<u>1.39</u>
<u>KGE</u>	<u>0.94</u>	<u>1.00</u>	<u>0.98</u>	<u>0.99</u>	<u>0.66</u>	<u>0.94</u>	<u>0.89</u>

410 *Author contributions.* DMB wrote the original draft and conducted the main analysis. RP contributed to the main analysis and, the writing  
411 and revision of manuscript. DMB, RP, VM, EJ, and RM conceptualized and designed the study. VM, EJ, EE, and RM contributed to the  
412 writing and revision of the manuscript. EE, AG, MD, JG, and JJ collected and provided in-situ data and offered expert feedback.

413 *Competing interests.* The authors declare that they have no conflict of interest.

414 *Acknowledgements.* This research was funded through "Horizon Europe funding program under Grant Agreement number 101081728"  
415 <https://doi.org/10.3030/101081728>, funded by the European Commission, as a part of the "Innovative tools to control organic matter and  
416 disinfection byproducts in drinking water" (intoDBP) project <https://intodbp.eu/>. [DMB and RF were funded by the Catalan Water Agency](#)  
417 [\(Agència Catalana de l'Aigua - ACA\) as a part of the NEREIDA project \(subvenció R+D+I de la convocatòria ACC/1362/2024\).](#)

## 418 References

- 419 Asadollah, S. B. H. S., Safaeinia, A., Jarahizadeh, S., Alcalá, F. J., Sharafati, A., and Jodar-Abellan, A.: Dissolved organic carbon estimation  
420 in lakes: Improving machine learning with data augmentation on fusion of multi-sensor remote sensing observations, *Water Research*,  
421 277, 123 350, 2025.
- 422 Bhateria, R. and Jain, D.: Water quality assessment of lake water: A review, *Sustainable Water Resources Management*, 2, 161–173,  
423 <https://doi.org/10.1007/s40899-015-0014-7>, 2016.
- 424 Biau, G. and Scornet, E.: A random forest guided tour, *Test*, 25, 197–227, 2016.
- 425 Breiman, L.: Random Forests, *Machine Learning*, 45, 5–32, <https://doi.org/10.1023/A:1010933404324>, 2001.
- 426 Chen, C., Chen, Q., Yao, S., He, M., Zhang, J., Li, G., and Lin, Y.: Combining physical-based model and machine  
427 learning to forecast chlorophyll-a concentration in freshwater lakes, *Science of The Total Environment*, 907, 168 097,  
428 <https://doi.org/10.1016/j.scitotenv.2023.168097>, 2024.
- 429 Chen, T. and Guestrin, C.: XGBoost: A Scalable Tree Boosting System, in: *Proceedings of the 22nd ACM SIGKDD International Conference*  
430 *on Knowledge Discovery and Data Mining*, pp. 785–794, <https://doi.org/10.1145/2939672.2939785>, 2016.
- 431 Cortes, C. and Vapnik, V.: Support-vector networks, *Machine Learning*, 20, 273–297, <https://doi.org/10.1007/BF00994018>, 1995.
- 432 Creed, I. F., Bergström, A.-K., Trick, C. G., Grimm, N. B., Hessen, D. O., Karlsson, J., Kidd, K. A., Kritzberg, E., McKnight, D. M.,  
433 Freeman, E. C., Senar, O. E., Andersson, A., Ask, J., Berggren, M., Cherif, M., Giesler, R., Hotchkiss, E. R., Kortelainen, P., Palta, M. M.,  
434 and Weyhenmeyer, G. A.: Global change-driven effects on dissolved organic matter composition: Implications for food webs of northern  
435 lakes, *Global Change Biology*, 24, 3692–3714, <https://doi.org/10.1111/gcb.14129>, 2018.
- 436 Downing, B. D., Pellerin, B. A., Bergamaschi, B. A., Saraceno, J. F., and Kraus, T. E. C.: Seeing the light: The effects of particles, dissolved  
437 materials, and temperature on in situ measurements of DOM fluorescence in rivers and streams, *Limnology and Oceanography: Methods*,  
438 10, 767–775, <https://doi.org/10.4319/lom.2012.10.767>, 2012.
- 439 Doyle, B. C., de Eyto, E., Dillane, M., Poole, R., McCarthy, V., Ryder, E., and Jennings, E.: Synchrony in catchment stream colour levels is  
440 driven by both local and regional climate, *Biogeosciences*, 16, 1053–1071, <https://doi.org/10.5194/bg-16-1053-2019>, 2019.
- 441 Duan, H., Cao, Z., Luo, J., and Shen, M.: AI-driven opportunities and challenges in lake remote sensing, *Information Geography*, p. 100014,  
442 <https://doi.org/10.1016/j.infgeo.2025.100014>, 2025.
- 443 Dubeux, J. C. B., Lira Junior, M. d. A., Simili, F. F., Bretas, I. L., Trumpp, K. R., Bizzuti, B. E., Garcia, L., Oduor, K. T., Queiroz, L. M. D.,  
444 Acuña, J. P., and Mendes, C. T. E.: Deep soil organic carbon: A review, *CABI Reviews*, 19, <https://doi.org/10.1079/cabireviews.2024.0024>,  
445 2024.
- 446 Fix, E.: *Discriminatory Analysis: Nonparametric Discrimination, Consistency Properties*, Tech. rep., USAF School of Aviation Medicine,  
447 1985.
- 448 Gobler, C. J.: Climate Change and Harmful Algal Blooms: Insights and perspective, *Harmful Algae*, 91, 101 731,  
449 <https://doi.org/10.1016/j.hal.2019.101731>, 2020.
- 450 Hanson, P. C., Stillman, A. B., Jia, X., Karpatne, A., Dugan, H. A., Carey, C. C., Stachelek, J., Ward, N. K., Zhang, Y., Read, J. S., and Kumar,  
451 V.: Predicting lake surface water phosphorus dynamics using process-guided machine learning, *Ecological Modelling*, 430, 109 136,  
452 <https://doi.org/10.1016/j.ecolmodel.2020.109136>, 2020.
- 453 Harkort, L. and Duan, Z.: Estimation of dissolved organic carbon from inland waters at a large scale using satellite data and machine learning  
454 methods, *Water Research*, 229, 119 478, <https://doi.org/10.1016/j.watres.2022.119478>, 2023.

455 Hersbach, H., Bell, B., Berrisford, P., Biavati, G., Horányi, A., Muñoz Sabater, J., Nicolas, J., Peubey, C., Radu, R., Rozum, I., Schep-  
456 ers, D., Simmons, A., Soci, C., Dee, D., and Thépaut, J.-N.: ERA5 hourly data on single levels from 1940 to present, Coperni-  
457 cus Climate Change Service (C3S) Climate Data Store (CDS), <https://doi.org/10.24381/cds.adbb2d47>, retrieved April 14, 2025, from  
458 <https://cds.climate.copernicus.eu/datasets/reanalysis-era5-single-levels?tab=overview>, Accessed on July 2025.

459 Herzsprung, P., Wentzky, V., Kamjunke, N., von Tümpling, W., Wilske, C., Friese, K., Boehrer, B., Reemtsma, T., Rinke, K., and Lechtenfeld,  
460 O. J.: Improved Understanding of Dissolved Organic Matter Processing in Freshwater Using Complementary Experimental and Machine  
461 Learning Approaches, *Environmental Science & Technology*, 54, 13 556–13 565, <https://doi.org/10.1021/acs.est.0c02383>, 2020.

462 Hipsey, M. R., Bruce, L. C., Boon, C., Busch, B., Carey, C. C., Hamilton, D. P., Hanson, P. C., Read, J. S., de Sousa, E., Weber, M.,  
463 and Winslow, L. A.: A General Lake Model (GLM 3.0) for linking with high-frequency sensor data from the Global Lake Ecological  
464 Observatory Network (GLEON), *Geoscientific Model Development*, 12, 473–523, <https://doi.org/10.5194/gmd-12-473-2019>, 2019.

465 Hollister, J. W., Milstead, W. B., and Kreakie, B. J.: Modeling lake trophic state: A random forest approach, *Ecosphere*, 7, e01 321,  
466 <https://doi.org/10.1002/ecs2.1321>, 2016.

467 Jennings, E., de Eyto, E., Moore, T., Dillane, M., Ryder, E., Allott, N., Nic Aonghusa, C., Rouen, M., Poole, R., and Pierson, D. C.: From  
468 Highs to Lows: Changes in Dissolved Organic Carbon in a Peatland Catchment and Lake Following Extreme Flow Events, *Water*, 12,  
469 2843, <https://doi.org/10.3390/w12102843>, 2020.

470 Kalbitz, K., Solinger, S., Park, J.-H., Michalzik, B., and Matzner, E.: Controls on the dynamics of dissolved organic matter in soils: A review,  
471 *Soil Science*, 165, 277–300, 2000.

472 Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., Ye, Q., and Liu, T.-Y.: LightGBM: A highly efficient gradient boosting decision  
473 tree, in: *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pp. 3149–3157, 2017.

474 Lake, P. S., Palmer, M. A., Biro, P., Cole, J., Covich, A. P., Dahm, C., Gibert, J., Goedkoop, W., Martens, K., and Verhoeven, J.: Global  
475 Change and the Biodiversity of Freshwater Ecosystems: Impacts on Linkages between Above-Sediment and Sediment Biota, *BioScience*,  
476 50, 1099–1107, [https://doi.org/10.1641/0006-3568\(2000\)050\[1099:GCATBO\]2.0.CO;2](https://doi.org/10.1641/0006-3568(2000)050[1099:GCATBO]2.0.CO;2), 2000.

477 Li, A., Zhao, X., Mao, R., Liu, H., and Qu, J.: Characterization of dissolved organic matter from surface waters with low to high  
478 dissolved organic carbon and the related disinfection byproduct formation potential, *Journal of Hazardous Materials*, 271, 228–235,  
479 <https://doi.org/10.1016/j.jhazmat.2014.02.009>, 2014.

480 Li, B., Yang, G., Wan, R., Dai, X., and Zhang, Y.: Comparison of random forests and other statistical methods for the prediction of lake water  
481 level: A case study of the Poyang Lake in China, *Hydrology Research*, 47, 69–83, <https://doi.org/10.2166/nh.2016.264>, 2016.

482 Li, M., del Giorgio, P. A., Parkes, A. H., and Prairie, Y. T.: The relative influence of topography and land cover on inorganic and or-  
483 ganic carbon exports from catchments in southern Quebec, Canada, *Journal of Geophysical Research: Biogeosciences*, 120, 2562–2578,  
484 <https://doi.org/10.1002/2015JG003073>, 2015.

485 Liu, D., Yu, S., Xiao, Q., Qi, T., and Duan, H.: Satellite estimation of dissolved organic carbon in eutrophic Lake Taihu, China, *Remote*  
486 *Sensing of Environment*, 264, 112 572, <https://doi.org/10.1016/j.rse.2021.112572>, 2021.

487 Marcé, R., Verdura, L., and Leung, N.: Dissolved organic matter spectroscopy reveals a hot spot of organic matter changes at the river-  
488 reservoir boundary, *Aquatic Sciences*, 83, 67, <https://doi.org/10.1007/s00027-021-00823-6>, 2021.

489 McCullough, I. M., Dugan, H. A., Farrell, K. J., Morales-Williams, A. M., Ouyang, Z., Roberts, D., Scordo, F., Bartlett, S. L., Burke, S. M.,  
490 Doubek, J. P., et al.: Dynamic modeling of organic carbon fates in lake ecosystems, *Ecological Modelling*, 386, 71–82, 2018.

491 Mercado-Bettín, D., Clayer, F., Shikhani, M., Moore, T. N., Frías, M. D., Jackson-Blake, L., Sample, J., Iturbide, M., Herrera, S., French,  
492 A. S., Norling, M. D., Rinke, K., and Marcé, R.: Forecasting water temperature in lakes and reservoirs using seasonal climate prediction,  
493 *Water Research*, 201, 117 286, <https://doi.org/10.1016/j.watres.2021.117286>, 2021.

494 Mi, C., Tilahun, A. B., Flörke, M., Dürr, H. H., and Rinke, K.: Climate warming effects in stratified reservoirs: Thorough assessment for  
495 opportunities and limits of machine learning techniques versus process-based models in thermal structure projections, *Journal of Cleaner*  
496 *Production*, 454, 142 347, <https://doi.org/10.1016/j.jclepro.2024.142347>, 2024.

497 Molnar, C.: *Interpretable machine learning*, Lulu. com, 2020.

498 Müller, M., D'Andrilli, J., Silverman, V., Bier, R. L., Barnard, M. A., Lee, M. C. M., Richard, F., Tanentzap, A. J., Wang, J., de Melo, M., and  
499 Lu, Y.: Machine-learning based approach to examine ecological processes influencing the diversity of riverine dissolved organic matter  
500 composition, *Frontiers in Water*, 6, <https://doi.org/10.3389/frwa.2024.1379284>, 2024.

501 Nearing, G. S., Kratzert, F., Sampson, A. K., Pelissier, C. S., Klotz, D., Frame, J. M., Prieto, C., and Gupta, H. V.: What  
502 Role Does Hydrological Science Play in the Age of Machine Learning?, *Water Resources Research*, 57, e2020WR028 091,  
503 <https://doi.org/10.1029/2020WR028091>, 2021.

504 Paíz, R., Pierson, D. C., Lindqvist, K., Naden, P. S., de Eyto, E., Dillane, M., McCarthy, V., Linnane, S., and Jennings, E.: Accounting for  
505 model parameter uncertainty provides more robust projections of dissolved organic carbon dynamics to aid drinking water management,  
506 *Water Research*, 276, 123 238, <https://doi.org/10.1016/j.watres.2025.123238>, 2025a.

507 Paíz, R., Thomas, R. Q., Carey, C. C., de Eyto, E., Jones, I. D., Delany, A. D., Poole, R., Nixon, P., Dillane, M., McCarthy, V., Linnane,  
508 S., and Jennings, E.: Near-term lake water temperature forecasts can be used to anticipate the ecological dynamics of freshwater species,  
509 *Ecosphere*, 16, e70 335, <https://doi.org/10.1002/ecs2.70335>, 2025b.

510 Prokhorenkova, L., Gusev, G., Vorobev, A., Dorogush, A. V., and Gulin, A.: CatBoost: Unbiased boosting with categorical features, *Tech.*  
511 *Rep.* arXiv:1706.09516, arXiv, <https://doi.org/10.48550/arXiv.1706.09516>, 2019.

512 Qi, Y.: *Random Forest for Bioinformatics*, pp. 307–323, Springer, [https://doi.org/10.1007/978-1-4419-9326-7\\_11](https://doi.org/10.1007/978-1-4419-9326-7_11), 2012.

513 Regier, P., Duggan, M., Myers-Pigg, A., and Ward, N.: Effects of random forest modeling decisions on biogeochemical time series predic-  
514 tions, *Limnology and Oceanography: Methods*, 21, 40–52, <https://doi.org/10.1002/lom3.10523>, 2023.

515 Rodriguez-Galiano, V., Sanchez-Castillo, M., Chica-Olmo, M., and Chica-Rivas, M.: Machine learning predictive models for mineral  
516 prospectivity: An evaluation of neural networks, random forest, regression trees and support vector machines, *Ore Geology Reviews*,  
517 71, 804–818, <https://doi.org/10.1016/j.oregeorev.2015.01.001>, 2015.

518 Rumpel, C. and Kögel-Knabner, I.: Deep soil organic matter—A key but poorly understood component of terrestrial C cycle, *Plant and Soil*,  
519 338, 143–158, <https://doi.org/10.1007/s11104-010-0391-5>, 2011.

520 Ryder, E., Jennings, E., de Eyto, E., Dillane, M., NicAonghusa, C., Pierson, D. C., Moore, K., Rouen, M., and Poole, R.: Temperature  
521 quenching of CDOM fluorescence sensors: Temporal and spatial variability in the temperature response and a recommended temperature  
522 correction equation, *Limnology and Oceanography: Methods*, 10, 1004–1010, <https://doi.org/10.4319/lom.2012.10.1004>, 2012.

523 Ryder, E., de Eyto, E., Dillane, M., Poole, R., and Jennings, E.: Identifying the role of environmental drivers in organic carbon export from  
524 a forested peat catchment, *Science of The Total Environment*, 490, 28–36, <https://doi.org/10.1016/j.scitotenv.2014.04.091>, 2014.

525 Solomon, C. T., Jones, S. E., Weidel, B., Buffam, I., Fork, M. L., Karlsson, J., Larsen, S., Lennon, J. T., Read, J. S., Sadro, S., and Saros,  
526 J. E.: Ecosystem consequences of changing inputs of terrestrial dissolved organic matter to lakes: Current knowledge and future challenges,  
527 *Ecosystems*, 18, 376–389, <https://doi.org/10.1007/s10021-015-9848-y>, 2015.

528 Sullivan, E.: Understanding from Machine Learning Models, *The British Journal for the Philosophy of Science*, 73, 109–133,  
529 <https://doi.org/10.1093/bjps/axz035>, 2022.

530 Toming, K., Kotta, J., Uuemaa, E., Sobek, S., Kutser, T., and Tranvik, L. J.: Predicting lake dissolved organic carbon at a global scale,  
531 *Scientific Reports*, 10, 8471, <https://doi.org/10.1038/s41598-020-65010-3>, 2020.

532 Šimek, K., Comerma, M., García, J.-C., Nedoma, J., Marcé, R., and Armengol, J.: The Effect of River Water Circulation on the Distri-  
533 bution and Functioning of Reservoir Microbial Communities as Determined by a Relative Distance Approach, *Ecosystems*, 14, 1–14,  
534 <https://doi.org/10.1007/s10021-010-9388-4>, 2011.

535 Watras, C., Hanson, P., Stacy, T., Morrison, K., Mather, J., Hu, Y.-H., and Milewski, P.: A temperature compensation method for CDOM  
536 fluorescence sensors in freshwater, *Limnology and Oceanography: Methods*, 9, 296–301, 2011.

537 Weyhenmeyer, G. A. and Karlsson, J.: Nonlinear response of dissolved organic carbon concentrations in boreal lakes to increasing tempera-  
538 tures, *Limnology and Oceanography*, 54, 2513–2519, [https://doi.org/10.4319/lo.2009.54.6\\_part\\_2.2513](https://doi.org/10.4319/lo.2009.54.6_part_2.2513), 2009.

539 Xenopoulos, M. A., Barnes, R. T., Boodoo, K. S., Butman, D., Catalán, N., D’Amario, S. C., Fasching, C., Kothawala, D. N., Pisani,  
540 O., Solomon, C. T., Spencer, R. G. M., Williams, C. J., and Wilson, H. F.: How humans alter dissolved organic matter composition  
541 in freshwater: Relevance for the Earth’s biogeochemistry, *Biogeochemistry*, 154, 323–348, <https://doi.org/10.1007/s10533-021-00753-3>,  
542 2021.

543 Zhang, D., Shi, K., Wang, W., Wang, X., Zhang, Y., Qin, B., Zhu, M., Dong, B., and Zhang, Y.: An optical mechanism-based  
544 deep learning approach for deriving water trophic state of China’s lakes from Landsat images, *Water Research*, 252, 121181,  
545 <https://doi.org/10.1016/j.watres.2024.121181>, 2024.

546 Zhang, Y., Yao, X., Wu, Q., Huang, Y., Zhou, Z., Yang, J., and Liu, X.: Turbidity prediction of lake-type raw water using random  
547 forest model based on meteorological data: A case study of Tai lake, China, *Journal of Environmental Management*, 290, 112657,  
548 <https://doi.org/10.1016/j.jenvman.2021.112657>, 2021.