



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

Daniel Mercado-Bettín¹, Ricardo Paíz^{2,3}, Valerie McCarthy², Eleanor Jennings³, Elvira de Eyto⁴, Angeles M. Gallegos⁵, Mary Dillanee⁴, Juan C. Garcia⁵, José J. Rodríguez⁵, and Rafael Marcé¹

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

²School of History and Geography, Dublin City University, D09 YT18, Dublin 9, Ireland

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

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

⁵Ens d'Abastament d'Aigua Ter-Llobregat, Ctra. Aigües, 6, 08440, Cardedeu, Spain

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

Abstract. Predicting water quality variables in lakes is critical for effective ecosystem management under climatic and human pressures. Dissolved organic matter (DOM) serves as an energy source for aquatic ecosystems and plays a key role in their biogeochemical cycles. However, predicting DOM is challenging due to complex interactions between multiple potential drivers in the aquatic environment and its surrounding terrestrial landscape. This study establishes an open and scalable workflow to identify potential drivers and predict fluorescent DOM (fDOM) in the surface layer of lakes by exploring the use of supervised machine learning models, including random forest, extreme gradient boosting, light gradient boosting, catboosting, k-nearest neighbors, support vector regression and linear model. It was validated in two contrasting systems: one natural lake in Ireland with a relatively undisturbed catchment, and one reservoir in Spain with a more human-influenced catchment. A total of 24 potential drivers were obtained from global reanalysis data, and lake and river process-based modelling. Partial dependence and SHapley Additive exPlanations (SHAP) analyses were conducted for the most influential drivers identified, with soil moisture, soil temperature, and Julian day being common to both study sites. The best prediction was found when using the CatBoost model (during hold-out testing period, Irish site: KGE > 0.69, r^2 > 0.51; Spanish site: KGE > 0.66, r^2 > 0.54). Interestingly, when only using drivers from globally accessible climate and soil reanalysis data, the prediction capacity was maintained at both sites, showcasing potential for scalability. Our findings highlight the complex interplay of environmental drivers and processes that govern DOM dynamics in lakes, and contribute to the modelling of carbon cycling in aquatic ecosystems.

1 Introduction

Lakes are an essential component of global biogeochemical cycles, sustain biodiversity, and provide critical ecosystem services, e.g., water supply, fishing and irrigation. However, their water quality is increasingly at risk due to climatic change and human pressures (Bhateria and Jain, 2016). A key water quality variable is dissolved organic matter (DOM), which influences light penetration, energy, oxygen dynamics and nutrient availability in any lake (Solomon et al., 2015). The dynamics of DOM in lakes are driven by both external processes in the terrestrial environment and internal processes. Land cover, climate, and



topography regulate carbon production in the catchment and carbon inputs into the lake (Li et al., 2015). In the water body, the quantity and quality of DOM are controlled by physical and biogeochemical mechanisms such as photodegradation, microbial processing and mixing dynamics, but can also be impacted by water abstraction or dam regulation (Xenopoulos et al., 2021).

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

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

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

This study introduces a workflow for predicting fluorescent DOM (fDOM) (a proxy for DOM) in lakes using a suite of supervised ML models driven by potential drivers either extracted from reanalysis data (climate and soil variables) or outputs from lake and catchment process-based models. The workflow was tested in two different study sites, one in Ireland and one in Spain, that represent contrasting settings for both the potential drivers and DOM dynamics. Model performance was first evaluated at each site using the most influential drivers to predict fDOM. Subsequently, a second simulation was performed using a subset of these drivers, limited to those sourced from reanalysis data, to evaluate the predictive capacity of the model in the context of higher scalability. The key research questions guiding this study were: (1) What are the most influential drivers of fDOM predictions, and how does their importance vary between two contrasting sites? (2) How accurately can supervised ML approaches predict lake fDOM driven by reanalysis-based data, and hydrologic and lake modelling outputs? (3) How easily can the workflow be reproduced and scaled to other sites?

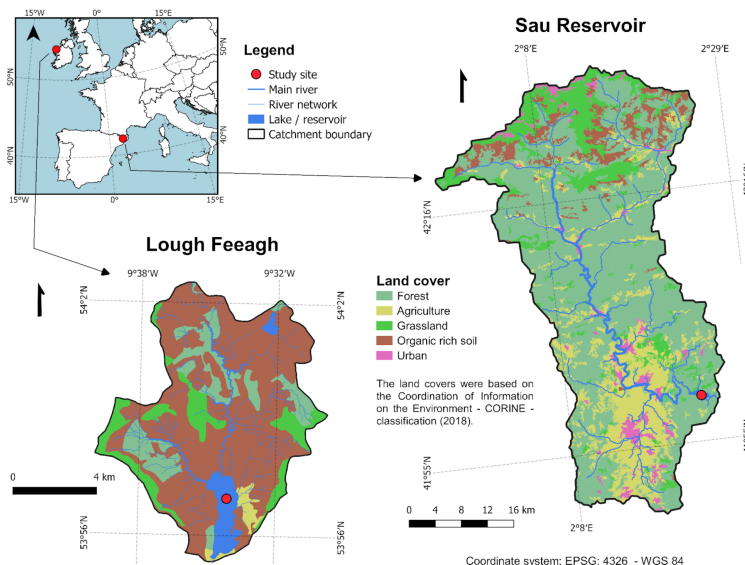


Figure 1. Two contrasting freshwater ecosystems. Lough Feeagh (Ireland) and Sau Reservoir (Spain) serve as contrasting sites for evaluating the predictive modelling of fDOM in lakes due to their distinct environmental and climatic conditions. The former is a humic and oligotrophic lake, dominated by a natural peatland catchment, and temperate oceanic climate, resulting in relatively higher levels of DOC during the whole year with a regular seasonality; the latter is an eutrophic heavily controlled reservoir, dominated by a highly anthropized 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: CORINE Land Cover 2018 <https://doi.org/10.2909/960998c1-1870-4e82-8051-6485205ebbac>

2 Materials and methods

2.1 Study sites

Lough Feeagh and Sau Reservoir are located in western Ireland (53° 56' N, 9° 35' W) and northeastern Spain (41° 58' N, 2° 23' E), respectively (Fig. 1). The study sites have contrasting attributes. Feeagh (depth of 46.8 m and area of 3.95 km²) is a monomictic and oligotrophic lake surrounded by a relatively undisturbed landscape, while Sau (depth of 70 m and area of 5.8 km²) is an eutrophic system subjected to human activities and water abstraction. Feeagh has two primary inflows, the Black and the Glenamong rivers, while Sau has one, the Ter river. The catchment of Feeagh is relatively small (84 km²), with mid-range hills, and dominated by peatland. The catchment of Sau, in contrast, is larger (1525 km²), with a varying topography and land uses (Fig.1; Fig. A1).

The dynamics of DOM in both study sites have been previously explored in Ryder et al. (2014) which identified soil temperature, river discharge and drought as important drivers in Feeagh, and in Marcé et al. (2021), which showed human activities were significant drivers in Sau. DOM in Feeagh is mostly driven by natural processes, while diffuse and point sources of or-



ganic matter and nutrients, e.g., wastewater effluents and agricultural runoff, are also important for Sau. Catchment hydrology is key for carbon transport into both study sites, and contributes to a distinct seasonality related to climate. Feeagh has a wet temperate climate, with cooler air temperatures and higher rainfall levels that occur on more than 75% of days in the year. The variability of carbon inputs reflects the sensitivity of a peatland-dominated landscape, which exacerbates climate-induced carbon release from the catchment into the lake. In contrast, Sau has a Mediterranean climate, characterised by hot, dry summers and mild, wet winters, dictating water availability, thermal stratification, and organic matter fluxes in the reservoir.

2.2 Prediction workflow

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

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

The same workflow was applied to both sites using the same data sources, allowing for comparison. Following the FAIR principles, all data and workflow scripts are available and fully reproducible in the following repository: https://github.com/danielmerbet/driver_attribution_fdom. Large language models were used in this study to optimise the codes, improve the final plots and, for basic proofreading of the text.

2.3 Data

2.3.1 Target variable (fDOM)

Daily surface fDOM values were computed from high-frequency data (2 minutes resolution) for both sites, for Feeagh measured at 0.9 m depth, and for Sau an average value was calculated between the depths of 0–5 m. fDOM data were expressed as quinine sulfate units (QSU) for the analysis. In Feeagh, the data spanned from 1st of May 2012 to 19th of November 2019 ($n = 2328$), and for Sau from 4th of February 2017 to 2nd of March 2020 ($n = 769$), with some gaps. All the other data (i.e., driver data) used in the workflow of this study were constrained by the availability of fDOM data.

In Feeagh, fDOM data were collected using a Seapoint UV fluorometer sensor (Seapoint Sensors Inc., Exeter, NH, USA) and water temperature data were measured using a Hach Environmental Hydrolab Data Sonde X5 (UK OTT Hydrometry Ltd). In Sau, fDOM and water temperature data were collected using a fDOM Digital Smart Sensor and Multiparameter Sonde

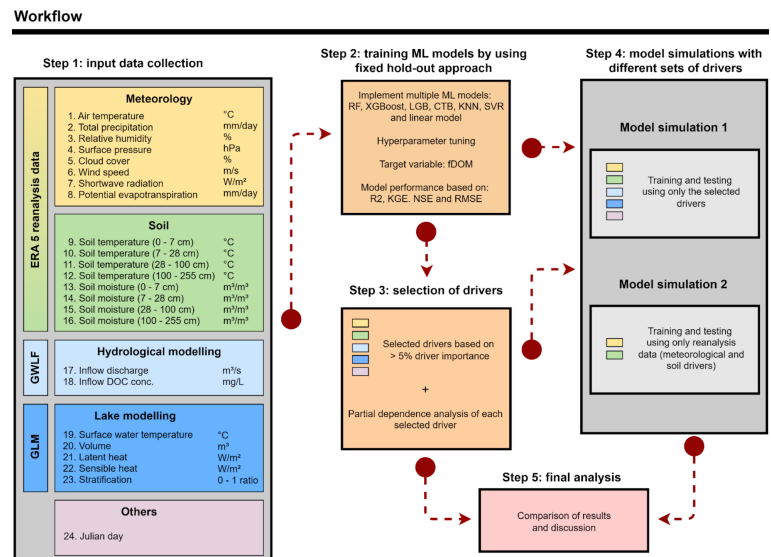


Figure 2. Workflow implemented for obtaining fDOM predictions in both study sites. The process consists of five steps: (1) Collecting input data representing all potential drivers of fDOM, including climate data (yellow) and soil data (green) from globally accessible reanalysis data, hydrologic model outputs (light blue), lake model outputs (blue), and other sources (light magenta). (2) Training the ML models by splitting the 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% of variable importance were retained for partial dependence analyses. (4) Running two simulations: using only the most influential selected drivers, and from these using only globally accessible reanalysis drivers and Julian day for ease and scalable implementation. (5) Analyzing and comparing the modelling results.

(YSI EXO sonde, Yellow Springs, OH, USA), respectively. Raw fDOM data were water temperature-corrected in both sites based on relationships established for each sensor (Ryder et al., 2012). Details about the fDOM corrections can be found in Supplementary Information and Figures A2 and A3.

2.3.2 Driver data

The input data for the workflow comprised 24 driving variables at each site. These were grouped into five categories: (1) meteorology, (2) soil, (3) process-based hydrological modelling, (4) process-based lake modelling and (5) Julian day. All input data variables, including their respective units and source, are displayed in Step 1 of Figure 2. Daily values of meteorology and soil variables were extracted from the European Centre for Medium-Range Weather Forecasts (ECMWF) Reanalysis v5 (ERA5) dataset (Hersbach et al., Accessed on July 2025). This gridded dataset provides (pseudo) observations at a global scale with a spatial resolution of 0.25°. Eight meteorological variables (traditionally employed in water modelling studies) and soil temperature and soil 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 each water body.



Daily values of inflow discharge and inflow DOC concentration into each site were generated using the Generalised Watershed Loading Functions Model (GWLf) coupled with a DOC module (GWLf-DOC). The GWLf-DOC version and calibration strategy applied are described in Paíz et al. (2025a). Calibration results can be found in the Supplementary Information. Daily values of five key lake variables (see Step 1, Fig. 2) were obtained from the General Lake Model (GLM) (Hipsey et al., 2019) run for both sites. Calibration strategies applied are described in Mercado-Bettín et al. (2021); Paíz et al. (2025b). Calibration results can be found in the Appendix. In addition, the cosine (to avoid an abrupt numerical change at every start of a year) of Julian day was included in the driver data inputs, given that seasonality is expected to influence DOM predictions.

2.4 Supervised Machine Learning

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

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

Seven ML and statistical models were used to predict fDOM dynamics: Random Forest (RF) (Breiman, 2001), eXtreme Gradient Boosting (XGBoost) (Chen and Guestrin, 2016), Light Gradient Boosting (LGB) (Ke et al., 2017), CatBoost (CTB) (Prokhorenkova et al., 2019), k-Nearest Neighbors (KNN) (Fix, 1985), Support Vector Regression (SVR) (Cortes and Vapnik, 1995) and linear model. RF, XGBoost, LGB and CTB can directly provide the importance of the features to predict the target variable, hence, only these four models were used to select the most important drivers to predict fDOM. For this, the increase in node purity was used for RF, and the gain contribution of each feature to the model for XGB, LGB and CTB. Further, hyperparameter tuning was implemented in all ML models to improve accuracy and generalisation. To extract the importance of the drivers, implement hyperparameter tuning, and predict fDOM, multiple R packages were used: caret, randomForest, xgboost, lightgbm, catboost, kkn and kernlab.



145 2.5 Partial dependence plots and SHAP analysis

146 To assess the influence of the most important drivers, we included partial dependence and SHapley Additive exPlanations
147 (SHAP) plots, using the Random Forest and CatBoost models, respectively. The partial dependence plots illustrate the indi-
148 vidual influence of each driver on fDOM predictions by varying the driver's values across its entire range while keeping all
149 other drivers constant in their average value. The SHAP plots measures how much a single driver (feature) value contributes
150 to moving the prediction away from the average value, the Y-axis has the input drivers ranked by overall importance (from top
151 to bottom), X-axis has the SHAP value representing the impact on model output for a single prediction, each point is a single
152 data instance and the color reflects the driver value (blue = low, red = high). To implement partial dependence plots and SHAP
153 analysis, the pdp package in R and the shap package in Python were used.

154 3 Results

155 3.1 Driver attribution

156 The most influential predictors of fDOM at each site were identified from all 24 potential drivers based on the 5% threshold
157 of the variable importance extracted from the RF, XGB, LGB and CTB models. This resulted in eight influential drivers being
158 identified for Feeagh and five for Sau (Figure 3), four of which were common to both sites.

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

165 3.2 Partial dependency on selected drivers

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

168 Seasonal patterns in fDOM concentrations were observed in both Feeagh and Sau, with higher values in winter and lower
169 in summer, as reflected in the influence of Julian day. However, the key predictors and their effects differed substantially,
170 shaped by contrasting catchment and climate characteristics. In Feeagh, where precipitation is relatively high and sustained
171 year-round, deep soil temperature (100–255 cm) was the dominant and potentially limiting predictor, with fDOM increasing
172 with temperature up to a threshold, beyond which a drop in the water table may counteract the effect. In addition, Feeagh
173 showed minimal influence of mid-depth soil temperature (28-100 cm) and solar radiation on fDOM. This temperature-fDOM
174 relationship was less relevant in Sau, where deep soil temperature (100-255 cm) had less explanatory power. Instead, fDOM

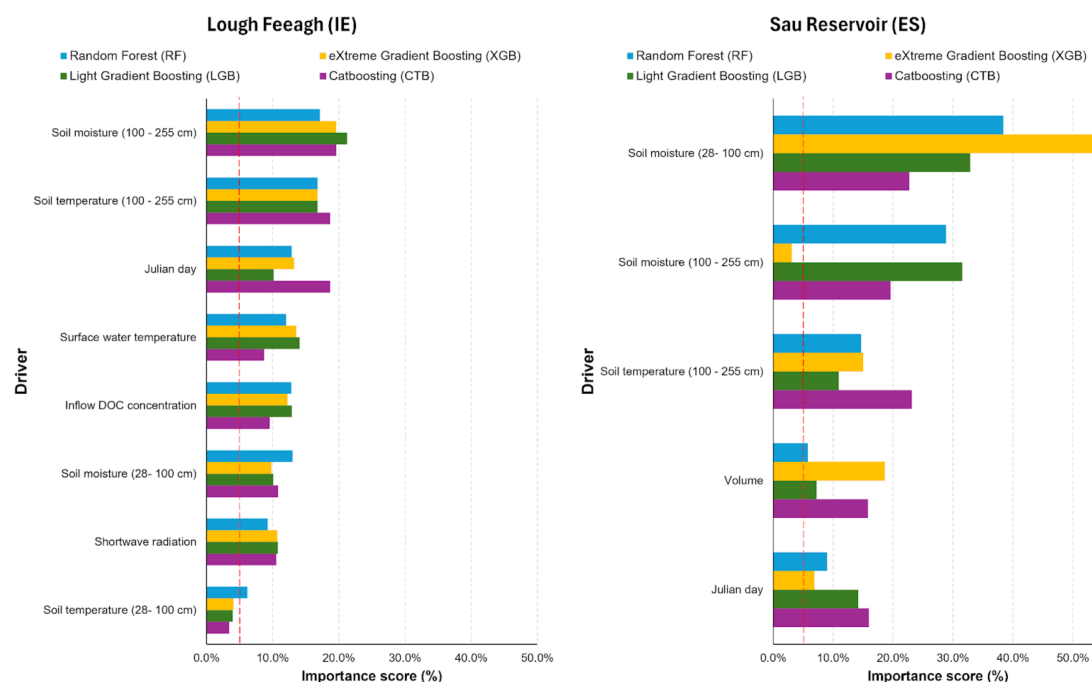


Figure 3. Selecting influential drivers to predict fDOM. 24 drivers from various sources were used to train the four ML models that directly provide feature importance: Random Forest (RF), eXtreme Gradient Boosting (XGBoost), Light Gradient Boosting (LGB), CatBoost (CTB). These included 8 climate variables, 8 soil variables, 2 outputs from hydrologic and water quality modelling, 5 outputs from lake modelling, and the cosine of the Julian day. The most relevant drivers were identified based on their contribution to node purity or gain contribution, with only those 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 Sau.

175 dynamics in Sau, were driven primarily by soil moisture at both 28-100 cm and 100-255 cm depths, potentially depicting a
 176 limiting condition by water stress.

177 Water availability also shaped the role of other predictors differently across the two sites. For instance, surface water tem-
 178 perature in Feeagh showed a clear threshold behaviour, with fDOM increasing relatively linearly beyond 6.5°C and stabilizing
 179 around 7.5°C, while in Sau, water volume acted as a surrogate for fDOM production. Lower volumes corresponded to reduced
 180 fDOM values, which increase and stabilize beyond 146 hm³. Interestingly, inflow DOC concentration influences Feeagh more
 181 than Sau, likely due to differences in hydroclimatological processes governing these relationships. These differences highlight
 182 how catchment water availability fundamentally alters the relative importance and behavior of fDOM drivers.

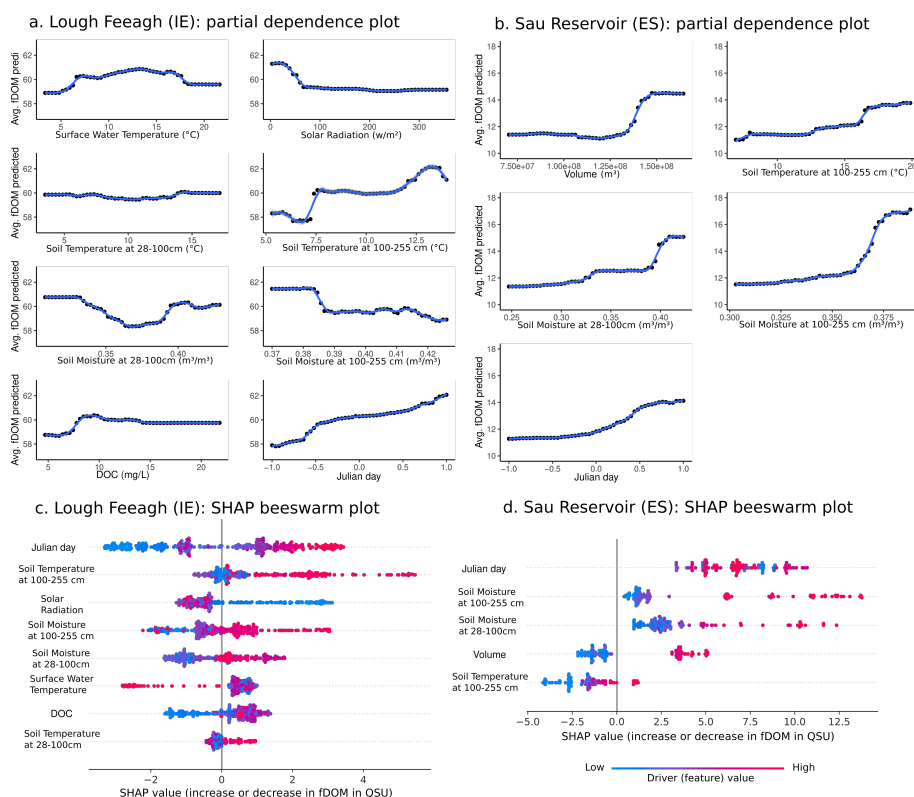


Figure 4. Driver influence in predicting fDOM. It presents the partial dependence plots for (a.) Lough Feeagh and (b.) Sau Reservoir based on the random forest (RF) model, and SHAP beeswarm plots for (c.) Lough Feeagh and (d.) Sau Reservoir based on the CatBoost (CTB) model. In Feeagh, soil temperature 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, Julian day was a relevant factor, although the influence of seasonality on fDOM predictions was more evident in Feeagh.

183 3.3 Predicting DOM using Supervised Machine Learning

184 Table 1 presents a comparative evaluation of the seven ML and statistical models employed in this study. In Lough Feeagh,
 185 CatBoost (CTB) demonstrated the best overall performance with the highest R^2 (0.51) and KGE (0.69), and a relatively low
 186 RMSE (8.29). Similarly, in Sau Reservoir, CTB again has the highest R^2 (0.54) and KGE (0.66), and one of the lowest RMSE
 187 (7.11). While some models like RF and LM showed moderate performance at Feeagh, others such as SVR and XGB performed
 188 poorly in Sau, with SVR even having a negative KGE (-0.82), suggesting significant model bias. Overall, CatBoost (CTB)
 189 consistently outperformed all other models across both sites, supporting its selection for fDOM prediction using the selected
 190 environmental drivers.

191 The results during the training phase (supplementary Table A1) confirm that most ML models, especially XGB, RF, and
 192 KNN, had a high performance during training. However, the performance dropped in some models (e.g., XGB at Sau Reser-



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, the best performance was found in the CatBoost model. 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	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

voir) during the testing phases, underscoring the importance of evaluating models on independent test data to assess general-
isability and overfitting. CatBoost (CTB), again, presented a more stable performance, showing slightly lower training metrics
(especially in Feeagh) compared with the other ML models and better generalisation when comparing with the metrics during
testing, supporting its suitability for prediction.

3.3.1 Model prediction using all drivers compared to a reduced set

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

For both lakes, the reduced driver models showed only a modest decline in predictive performance during the testing period.
For example, in Feeagh, the model using all influential drivers achieved $R^2 = 0.51$ and $KGE = 0.69$, while the reduced model
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
model maintained a comparable $R^2 = 0.50$ and $KGE = 0.65$. Although the training performance was higher in Sau compared
with the testing performance, indicating potential overfitting, the CatBoost model provided informative and generalisable
predictions in both study sites.

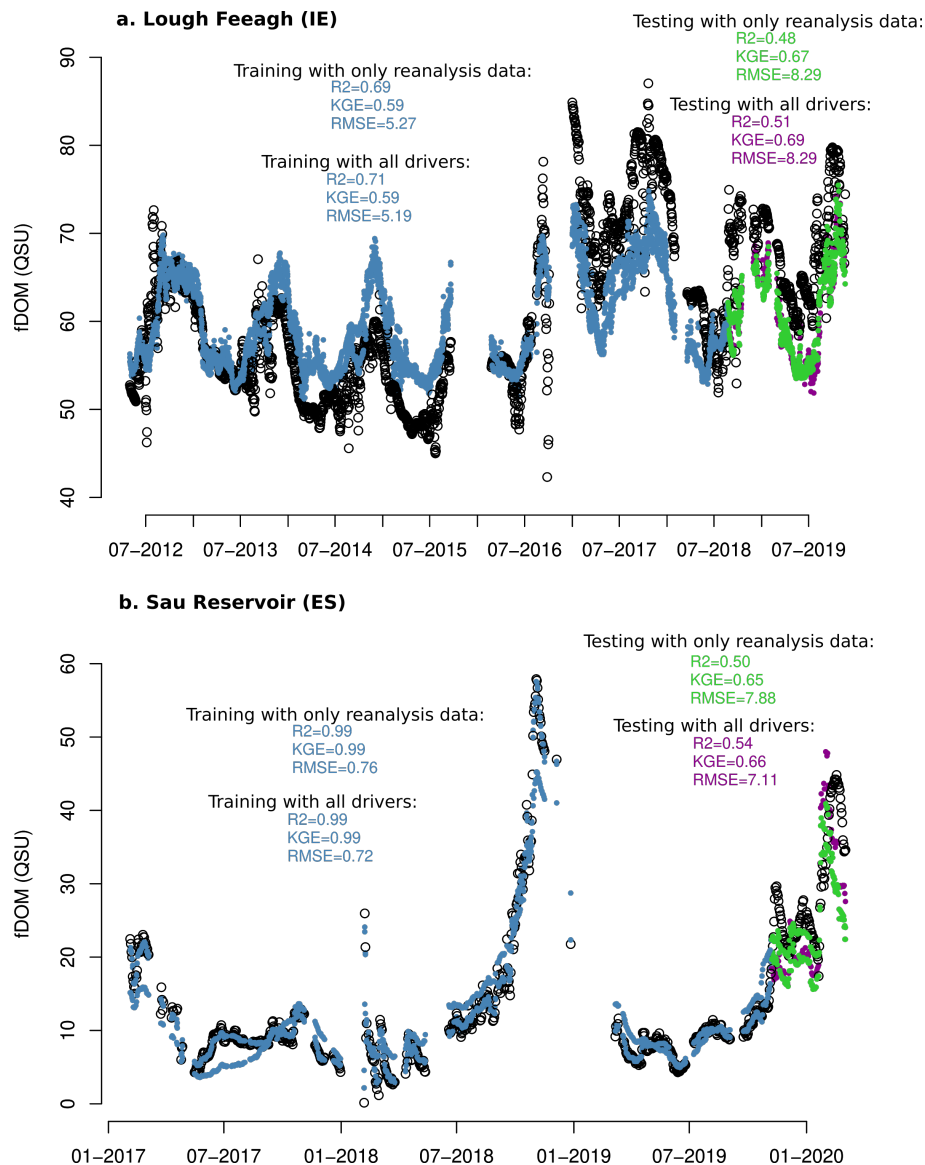


Figure 5. fDOM predictions obtained from the CatBoost model using different driver sets. Here the training (blue, 85% of the time series) and testing (violet or green, 15% of the time series) periods are shown. The training and testing periods are completely independent and follow the hold-out period method, and model performance metrics (R^2 , KGE, and RMSE) are calculated for both periods. a. Training and testing results for Feeagh: using the 8 most influential drivers (violet) and using only easily accessible and reanalysis data (soil temperature, soil moisture and Julian day) from those influential drivers (green); b. Training and testing results for Sau: using the 5 most influential drivers (violet) and using only easily accessible and reanalysis data (soil temperature, soil moisture and Julian day) from those influential drivers (green).



209 4 Discussion

210 4.1 Driver attribution in contrasting conditions

211 The most influential drivers for fDOM identified for each site suggest potential carbon-producing processes that drive DOM
 212 dynamics in each lake. Interestingly, the most influential drivers were related to temperature and moisture at the deepest
 213 soil layers, indicating a common relevance for carbon inputs from terrestrial processes in both catchments. The Irish site is
 214 dominated by peat, a highly organic soil known to be sensitive to temperature-induced carbon release, especially as temperature
 215 levels rise. This is supported by the relation of soil temperatures with fDOM at this site (Ryder et al., 2014). In contrast, water
 216 availability constraints are much more pronounced in drier soils such as those of the Spanish site, as is validated by the partial
 217 dependence relationship of soil moisture with fDOM (Šimek et al., 2011). A plausible explanation is that, in Sau, increased soil
 218 moisture could boost biological activity, enhancing fDOM production, while in Feeagh, soil moisture showed a bell-shaped
 219 relationship, suggesting a more nuanced interplay between oxygen availability and microbial processes (Fig 4).

220 DOM dynamics are driven by physical and biogeochemical processes in the soil that are sensitive to changes in temperature
 221 and moisture, e.g., microbial processes that break down organic matter (Kalbitz et al., 2000). The fact that the deepest layers
 222 of the soil were more important in the model for both sites than those shallower could be linked to potential carbon attenuation
 223 processes, such as soil organic matter decomposition and retention in the soil, including sorption processes (Dubeux et al.,
 224 2024; Rumpel and Kögel-Knabner, 2011). In any case, the production of carbon in the catchment that eventually ends up in
 225 the lake requires concurrent downstream transport, governed by rainfall events.

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

234 Seasonality plays a crucial role in fDOM predictions, as evidenced by the relationship of the Julian day driver with DOM
 235 dynamics at both sites (Fig. 4). At the Irish site, DOM seasonality is primarily shaped by natural environmental processes,
 236 whereas in the Spanish site, human influence plays a much greater role. This distinction helps explain why surface lake water
 237 temperature and solar radiation, two variables typically linked to strong seasonal patterns, were important only for Feeagh,
 238 while reservoir volume was significant only for Sau (Fig. 3 and 4). Volume and soil variables produce a similar effect on fDOM
 239 as Julian day at Sau, given that higher volumes closer to the winter season can lead to higher fDOM values. Incorporating Julian
 240 day into the workflow offers a simple yet effective way to represent seasonality, potentially replacing seasonal variables (e.g.,
 241 air temperature) (see the correlation matrix of all drivers in supplementary Fig. A5). This proves that the use of machine
 242 learning approaches opens up opportunities to assess diverse drivers under contrasting conditions.



Improving the accuracy of DOM predictions in lakes can enhance efforts to reduce further water quality deterioration and support lake management. This study demonstrated a feasible approach for simulating daily fDOM in two contrasting lakes, especially when using the Catboost model given its good generalisability. The performance metrics (Fig. 5) obtained at each site for the different model simulations lie in a similar or better range than comparable studies that modelled carbon dynamics in lakes (e.g., Harkort and Duan, 2023; Liu et al., 2021; Zhang et al., 2021, 2024). It is important to be aware, however, that previous studies are based on different frameworks. These include variations in the machine learning algorithms used, the target variable for quantifying carbon dynamics, with most studies having focused on DOC, whereas fDOM is the target variable here, as well as differences in input driver data and site-specific conditions.

4.2 Scalability

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

4.3 Limitations and future research

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

The workflow presented here is not recommended for climate change studies, as the drivers of DOM variability can significantly change under entirely new and unrecorded climatic conditions. Consequently, supervised machine learning may fail to capture the signal from the time series. Moreover, the method's reliance on historical patterns limits its ability to extrapolate



beyond the range of observed environmental conditions (Mi et al., 2024). Future research could expand the application of this framework to a broader range of lakes, integrate additional drivers such as remote sensing-derived terrestrial and aquatic quantity and quality parameters (Duan et al., 2025).

5 Conclusions

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

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

Code and data availability. All data and codes used in this study are available in this repository: https://github.com/danielmerbet/driver_attribution_fdom. A full and detailed README file and DOI link will be provided after the review process.

Appendix A: Supplementary information

Topography of Feeagh and Sau catchments

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

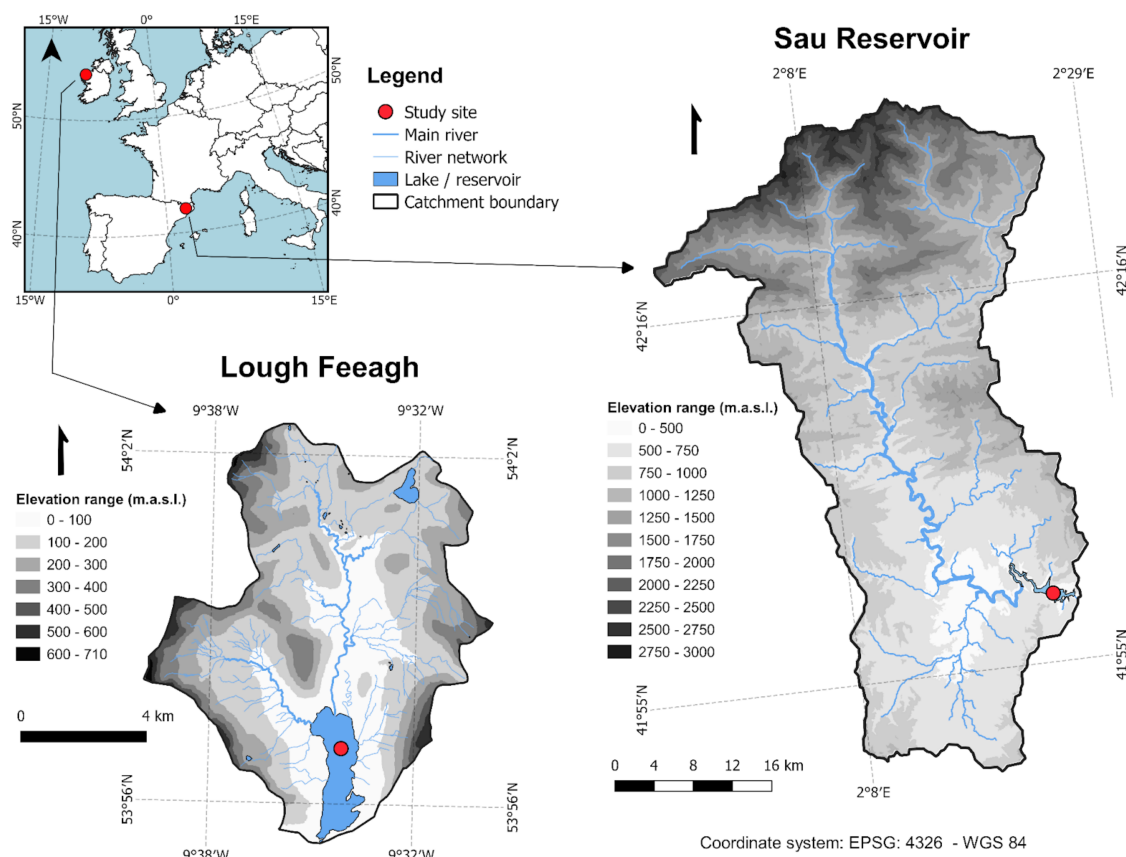


Figure A1. Elevation range for the two contrasting freshwater ecosystems: Lough Feeagh and Sau Reservoir.

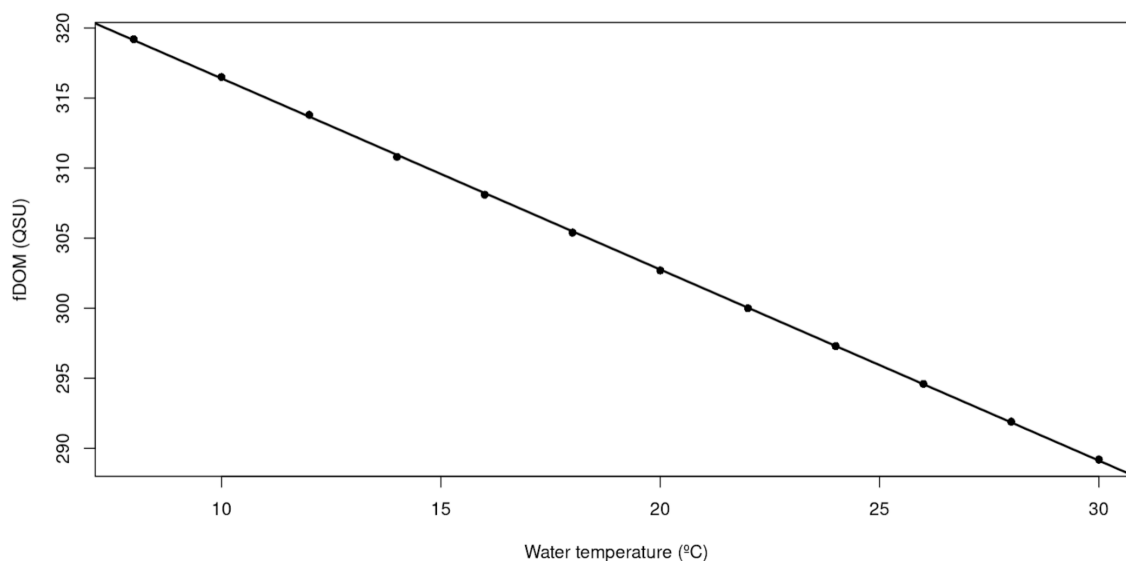


Figure A2. fDOM relation with temperature with a sample of 300 QSU provided by the manufacturer of the fDOM sensor in Sau.

295 fDOM correction for both study sites

296 In Feeagh, fDOM data was corrected for the temperature quenching effect in previous scientific studies (Doyle et al., 2019;
 297 Ryder et al., 2012).

298 In Sau, the fDOM data were corrected for the temperature quenching effect, following a test provided by the fDOM sensor
 299 manufacturer, where they use a 300 QSU sample of water and change the temperature to get the effect of temperature in the
 300 measurement, results can be found in Figure A2.

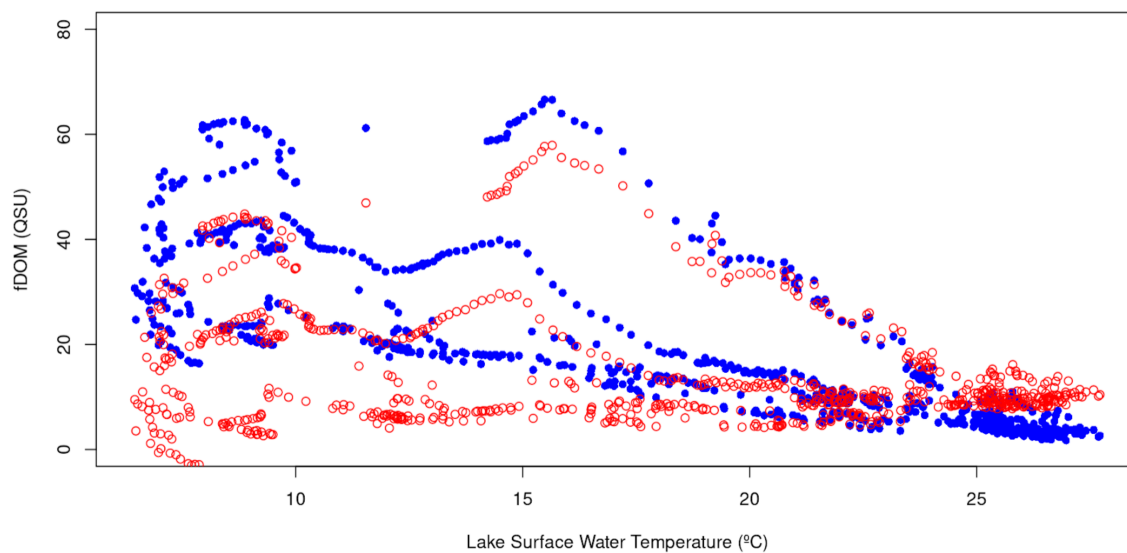


Figure A3. Uncorrected (blue) and corrected (red) fDOM data for Sau.

301 Then, fDOM data in Sau was corrected following this linear regression and surface water temperature on the lake. Figure
 302 A3 presents the uncorrected values in blue and corrected values in red (8 negative values were removed from the total sample
 303 of 777)



304 **Hydrologic Modelling**

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

310 GWLF-DOC was applied to Feeagh based on previous model applications in the Irish catchment (Paíz et al., 2025a), for
 311 which measured discharge data were used to calibrate and validate the hydrology (2013-2018 and 2019-2023, respectively),
 312 and DOC concentration data were used to calibrate the DOC module (2016-2023). In Sau, observed inflow discharge data
 313 were used to calibrate and validate the hydrology (2008-2011 and 2011-2024, respectively), and measured DOC concentration
 314 data were used to calibrate and validate the DOC module (2008-2014 and 2016-2018, respectively) using the same calibration
 315 strategy than for the Irish site. Calibration results were satisfactory for both hydrology (Feeagh: $R^2 = 0.64$ and $NSE = 0.64$;
 316 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,
 317 validation results were satisfactory for both hydrology (Feeagh: $R^2 = 0.60$ and $NSE = 0.60$; Sau: $R^2 = 0.42$ and $NSE = 0.42$)
 318 and DOC in the case of Sau ($R^2 = 0.50$ and $NSE = 0.46$).

319 **Lake Modelling**

320 Daily time series of 5 key lake variables (see Fig. 2) were obtained from the General Lake Model (GLM) run for each site. GLM
 321 is an open-source, one-dimensional hydrodynamic model designed to simulate the vertical stratification and water balance of
 322 lakes and reservoirs. It calculates vertical profiles of temperature, and density by accounting for factors such as inflows and
 323 outflows, mixing processes, and surface heating and cooling (Hipsey et al., 2019). GLM was calibrated and validated by
 324 evaluating the fit of modelled water temperature against measured water temperature profile data in Feeagh (2010-2015 and
 325 2016-2017, respectively) and Sau (1997-2007 and 2008-2018, respectively). The calibration strategy was based on previous
 326 lake modelling deployments at each site (Mercado-Bettín et al., 2021; Paíz et al., 2025a). Model performance was satisfactory
 327 for both sites.

328 **Comparison of validation methods using Random Forest**

329 To pick the most suitable validation method, we implemented hold-out period method used in the main manuscript, k-fold
 330 cross-validation method using $k=5$, and rolling window cross-validation using training size of two year, testing size of one
 331 year, and a shift window every 90 days. Results are shown in Figure A4.

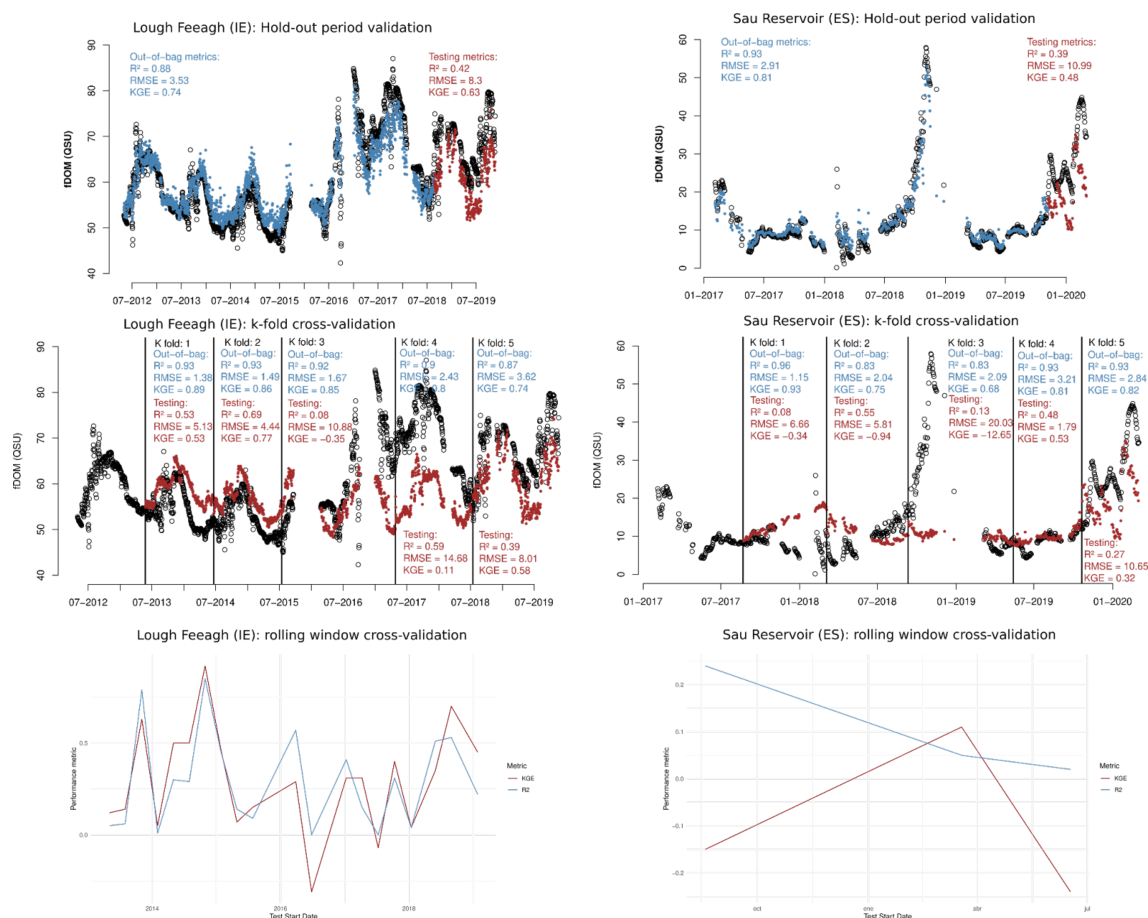


Figure A4. Comparison of validation methods using random forest for both study sites: hold-out period method used in the main manuscript, k-fold cross-validation method using $k=5$, and rolling window cross-validation using training size of two year, testing size of one year, and a shift window every 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.



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

Correlation matrix of all drivers is presented in Figure A5.

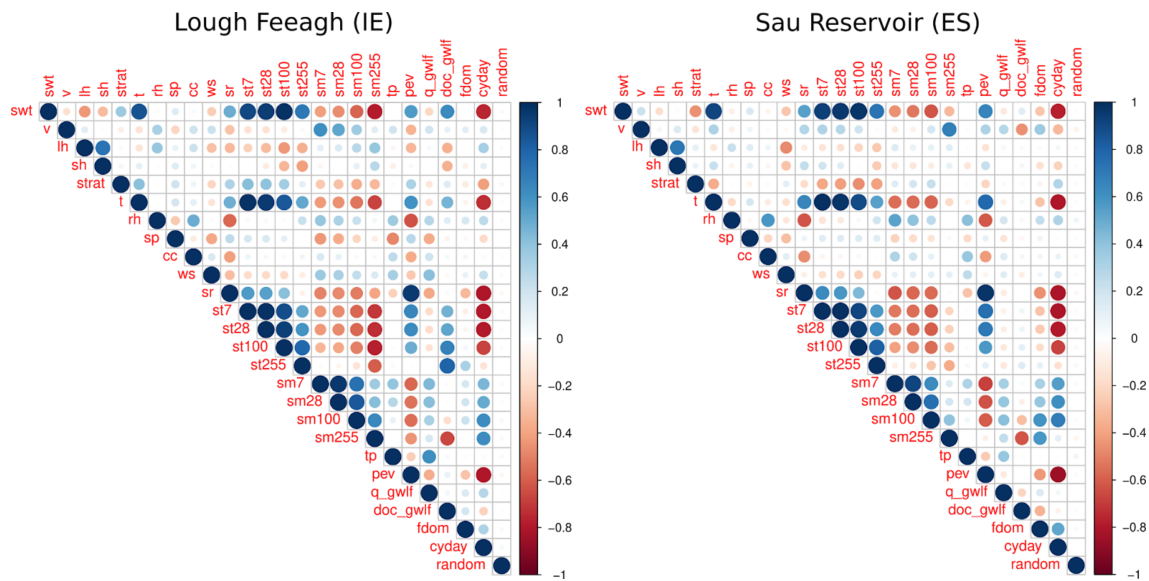


Figure A5. Correlation matrix of all drivers and fDOM for Feeagh and Sau.

333



334 **Model comparison during training to predict fDOM using all selected drivers in both study sites**

335 Resulting metrics during training period for all models are shown in Table A1.

Table A1. 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	0.88	0.71	0.22	0.99	0.89	0.99	1.00	0.98	0.99	0.67	0.98	0.98
RMSE	0.97	0.40	3.35	5.19	7.87	0.92	2.95	0.81	0.11	1.48	0.72	5.43	1.17	1.38
KGE	0.95	0.99	0.79	0.59	0.25	0.98	0.90	0.97	1.00	0.97	0.99	0.75	0.99	0.98



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

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

340 *Acknowledgements.* This research was funded through "Horizon Europe funding program under Grant Agreement number 101081728"
341 <https://doi.org/10.3030/101081728>, funded by the European Commission, as a part of the "Innovative tools to control organic matter and
342 disinfection byproducts in drinking water" (intoDBP) project <https://intodbp.eu/>



343 References

- 344 Asadollah, S. B. H. S., Safaeinia, A., Jarahizadeh, S., Alcalá, F. J., Sharafati, A., and Jodar-Abellan, A.: Dissolved organic carbon estimation
 345 in lakes: Improving machine learning with data augmentation on fusion of multi-sensor remote sensing observations, *Water Research*,
 346 277, 123 350, 2025.
- 347 Bhateria, R. and Jain, D.: Water quality assessment of lake water: A review, *Sustainable Water Resources Management*, 2, 161–173,
 348 <https://doi.org/10.1007/s40899-015-0014-7>, 2016.
- 349 Biau, G. and Scornet, E.: A random forest guided tour, *Test*, 25, 197–227, 2016.
- 350 Breiman, L.: Random Forests, *Machine Learning*, 45, 5–32, <https://doi.org/10.1023/A:1010933404324>, 2001.
- 351 Chen, C., Chen, Q., Yao, S., He, M., Zhang, J., Li, G., and Lin, Y.: Combining physical-based model and machine
 352 learning to forecast chlorophyll-a concentration in freshwater lakes, *Science of The Total Environment*, 907, 168 097,
 353 <https://doi.org/10.1016/j.scitotenv.2023.168097>, 2024.
- 354 Chen, T. and Guestrin, C.: XGBoost: A Scalable Tree Boosting System, in: *Proceedings of the 22nd ACM SIGKDD International Conference*
 355 *on Knowledge Discovery and Data Mining*, pp. 785–794, <https://doi.org/10.1145/2939672.2939785>, 2016.
- 356 Cortes, C. and Vapnik, V.: Support-vector networks, *Machine Learning*, 20, 273–297, <https://doi.org/10.1007/BF00994018>, 1995.
- 357 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.,
 358 Freeman, E. C., Senar, O. E., Andersson, A., Ask, J., Berggren, M., Cherif, M., Giesler, R., Hotchkiss, E. R., Kortelainen, P., Palta, M. M.,
 359 and Weyhenmeyer, G. A.: Global change-driven effects on dissolved organic matter composition: Implications for food webs of northern
 360 lakes, *Global Change Biology*, 24, 3692–3714, <https://doi.org/10.1111/gcb.14129>, 2018.
- 361 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
 362 materials, and temperature on in situ measurements of DOM fluorescence in rivers and streams, *Limnology and Oceanography: Methods*,
 363 10, 767–775, <https://doi.org/10.4319/lom.2012.10.767>, 2012.
- 364 Doyle, B. C., de Eyto, E., Dillane, M., Poole, R., McCarthy, V., Ryder, E., and Jennings, E.: Synchrony in catchment stream colour levels is
 365 driven by both local and regional climate, *Biogeosciences*, 16, 1053–1071, <https://doi.org/10.5194/bg-16-1053-2019>, 2019.
- 366 Duan, H., Cao, Z., Luo, J., and Shen, M.: AI-driven opportunities and challenges in lake remote sensing, *Information Geography*, p. 100014,
 367 <https://doi.org/10.1016/j.infgeo.2025.100014>, 2025.
- 368 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.,
 369 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>,
 370 2024.
- 371 Fix, E.: Discriminatory Analysis: Nonparametric Discrimination, Consistency Properties, Tech. rep., USAF School of Aviation Medicine,
 372 1985.
- 373 Gobler, C. J.: Climate Change and Harmful Algal Blooms: Insights and perspective, *Harmful Algae*, 91, 101 731,
 374 <https://doi.org/10.1016/j.hal.2019.101731>, 2020.
- 375 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,
 376 V.: Predicting lake surface water phosphorus dynamics using process-guided machine learning, *Ecological Modelling*, 430, 109 136,
 377 <https://doi.org/10.1016/j.ecolmodel.2020.109136>, 2020.
- 378 Harkort, L. and Duan, Z.: Estimation of dissolved organic carbon from inland waters at a large scale using satellite data and machine learning
 379 methods, *Water Research*, 229, 119 478, <https://doi.org/10.1016/j.watres.2022.119478>, 2023.



- 380 Hersbach, H., Bell, B., Berrisford, P., Biavati, G., Horányi, A., Muñoz Sabater, J., Nicolas, J., Peubey, C., Radu, R., Rozum, I., Schep-
 381 ers, D., Simmons, A., Soci, C., Dee, D., and Thépaut, J.-N.: ERA5 hourly data on single levels from 1940 to present, Coperni-
 382 cus Climate Change Service (C3S) Climate Data Store (CDS), <https://doi.org/10.24381/cds.adbb2d47>, retrieved April 14, 2025, from
 383 <https://cds.climate.copernicus.eu/datasets/reanalysis-era5-single-levels?tab=overview>, Accessed on July 2025.
- 384 Herzsprung, P., Wentzky, V., Kamjunke, N., von Tümpling, W., Wilske, C., Friese, K., Boehrer, B., Reemtsma, T., Rinke, K., and Lechtenfeld,
 385 O. J.: Improved Understanding of Dissolved Organic Matter Processing in Freshwater Using Complementary Experimental and Machine
 386 Learning Approaches, *Environmental Science & Technology*, 54, 13 556–13 565, <https://doi.org/10.1021/acs.est.0c02383>, 2020.
- 387 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.,
 388 and Winslow, L. A.: A General Lake Model (GLM 3.0) for linking with high-frequency sensor data from the Global Lake Ecological
 389 Observatory Network (GLEON), *Geoscientific Model Development*, 12, 473–523, <https://doi.org/10.5194/gmd-12-473-2019>, 2019.
- 390 Hollister, J. W., Milstead, W. B., and Kreakie, B. J.: Modeling lake trophic state: A random forest approach, *Ecosphere*, 7, e01 321,
 391 <https://doi.org/10.1002/ecs2.1321>, 2016.
- 392 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
 393 Highs to Lows: Changes in Dissolved Organic Carbon in a Peatland Catchment and Lake Following Extreme Flow Events, *Water*, 12,
 394 2843, <https://doi.org/10.3390/w12102843>, 2020.
- 395 Kalbitz, K., Solinger, S., Park, J.-H., Michalzik, B., and Matzner, E.: Controls on the dynamics of dissolved organic matter in soils: A review,
 396 *Soil Science*, 165, 277–300, 2000.
- 397 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
 398 tree, in: *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pp. 3149–3157, 2017.
- 399 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
 400 Change and the Biodiversity of Freshwater Ecosystems: Impacts on Linkages between Above-Sediment and Sediment Biota, *BioScience*,
 401 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.
- 402 Li, A., Zhao, X., Mao, R., Liu, H., and Qu, J.: Characterization of dissolved organic matter from surface waters with low to high
 403 dissolved organic carbon and the related disinfection byproduct formation potential, *Journal of Hazardous Materials*, 271, 228–235,
 404 <https://doi.org/10.1016/j.jhazmat.2014.02.009>, 2014.
- 405 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
 406 level: A case study of the Poyang Lake in China, *Hydrology Research*, 47, 69–83, <https://doi.org/10.2166/nh.2016.264>, 2016.
- 407 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-
 408 ganic carbon exports from catchments in southern Quebec, Canada, *Journal of Geophysical Research: Biogeosciences*, 120, 2562–2578,
 409 <https://doi.org/10.1002/2015JG003073>, 2015.
- 410 Liu, D., Yu, S., Xiao, Q., Qi, T., and Duan, H.: Satellite estimation of dissolved organic carbon in eutrophic Lake Taihu, China, *Remote*
 411 *Sensing of Environment*, 264, 112 572, <https://doi.org/10.1016/j.rse.2021.112572>, 2021.
- 412 Marcé, R., Verdura, L., and Leung, N.: Dissolved organic matter spectroscopy reveals a hot spot of organic matter changes at the river–
 413 reservoir boundary, *Aquatic Sciences*, 83, 67, <https://doi.org/10.1007/s00027-021-00823-6>, 2021.
- 414 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.,
 415 Doubek, J. P., et al.: Dynamic modeling of organic carbon fates in lake ecosystems, *Ecological Modelling*, 386, 71–82, 2018.



- 416 Mercado-Bettín, D., Clayer, F., Shikhaní, M., Moore, T. N., Frías, M. D., Jackson-Blake, L., Sample, J., Iturbide, M., Herrera, S., French,
417 A. S., Norling, M. D., Rinke, K., and Marcé, R.: Forecasting water temperature in lakes and reservoirs using seasonal climate prediction,
418 Water Research, 201, 117 286, <https://doi.org/10.1016/j.watres.2021.117286>, 2021.
- 419 Mi, C., Tilahun, A. B., Flörke, M., Dürr, H. H., and Rinke, K.: Climate warming effects in stratified reservoirs: Thorough assessment for
420 opportunities and limits of machine learning techniques versus process-based models in thermal structure projections, Journal of Cleaner
421 Production, 454, 142 347, <https://doi.org/10.1016/j.jclepro.2024.142347>, 2024.
- 422 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
423 Lu, Y.: Machine-learning based approach to examine ecological processes influencing the diversity of riverine dissolved organic matter
424 composition, Frontiers in Water, 6, <https://doi.org/10.3389/frwa.2024.1379284>, 2024.
- 425 Nearing, G. S., Kratzert, F., Sampson, A. K., Pelissier, C. S., Klotz, D., Frame, J. M., Prieto, C., and Gupta, H. V.: What
426 Role Does Hydrological Science Play in the Age of Machine Learning?, Water Resources Research, 57, e2020WR028091,
427 <https://doi.org/10.1029/2020WR028091>, 2021.
- 428 Paíz, R., Pierson, D. C., Lindqvist, K., Naden, P. S., de Eyto, E., Dillane, M., McCarthy, V., Linnane, S., and Jennings, E.: Accounting for
429 model parameter uncertainty provides more robust projections of dissolved organic carbon dynamics to aid drinking water management,
430 Water Research, 276, 123 238, <https://doi.org/10.1016/j.watres.2025.123238>, 2025a.
- 431 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,
432 S., and Jennings, E.: Near-term lake water temperature forecasts can be used to anticipate the ecological dynamics of freshwater species,
433 Ecosphere, 16, e70 335, <https://doi.org/10.1002/ecs2.70335>, 2025b.
- 434 Prokhorenkova, L., Gusev, G., Vorobev, A., Dorogush, A. V., and Gulin, A.: CatBoost: Unbiased boosting with categorical features, Tech.
435 Rep. arXiv:1706.09516, arXiv, <https://doi.org/10.48550/arXiv.1706.09516>, 2019.
- 436 Qi, Y.: Random Forest for Bioinformatics, pp. 307–323, Springer, https://doi.org/10.1007/978-1-4419-9326-7_11, 2012.
- 437 Regier, P., Duggan, M., Myers-Pigg, A., and Ward, N.: Effects of random forest modeling decisions on biogeochemical time series predic-
438 tions, Limnology and Oceanography: Methods, 21, 40–52, <https://doi.org/10.1002/lom3.10523>, 2023.
- 439 Rodriguez-Galiano, V., Sanchez-Castillo, M., Chica-Olmo, M., and Chica-Rivas, M.: Machine learning predictive models for mineral
440 prospectivity: An evaluation of neural networks, random forest, regression trees and support vector machines, Ore Geology Reviews,
441 71, 804–818, <https://doi.org/10.1016/j.oregeorev.2015.01.001>, 2015.
- 442 Rumpel, C. and Kögel-Knabner, I.: Deep soil organic matter—A key but poorly understood component of terrestrial C cycle, Plant and Soil,
443 338, 143–158, <https://doi.org/10.1007/s11104-010-0391-5>, 2011.
- 444 Ryder, E., Jennings, E., de Eyto, E., Dillane, M., NicAonghusa, C., Pierson, D. C., Moore, K., Rouen, M., and Poole, R.: Temperature
445 quenching of CDOM fluorescence sensors: Temporal and spatial variability in the temperature response and a recommended temperature
446 correction equation, Limnology and Oceanography: Methods, 10, 1004–1010, <https://doi.org/10.4319/lom.2012.10.1004>, 2012.
- 447 Ryder, E., de Eyto, E., Dillane, M., Poole, R., and Jennings, E.: Identifying the role of environmental drivers in organic carbon export from
448 a forested peat catchment, Science of The Total Environment, 490, 28–36, <https://doi.org/10.1016/j.scitotenv.2014.04.091>, 2014.
- 449 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,
450 J. E.: Ecosystem consequences of changing inputs of terrestrial dissolved organic matter to lakes: Current knowledge and future challenges,
451 Ecosystems, 18, 376–389, <https://doi.org/10.1007/s10021-015-9848-y>, 2015.
- 452 Sullivan, E.: Understanding from Machine Learning Models, The British Journal for the Philosophy of Science, 73, 109–133,
453 <https://doi.org/10.1093/bjps/axz035>, 2022.



- 454 Toming, K., Kotta, J., Uemaa, E., Sobek, S., Kutser, T., and Tranvik, L. J.: Predicting lake dissolved organic carbon at a global scale,
 455 Scientific Reports, 10, 8471, <https://doi.org/10.1038/s41598-020-65010-3>, 2020.
- 456 Šimek, K., Comerma, M., García, J.-C., Nedoma, J., Marcé, R., and Armengol, J.: The Effect of River Water Circulation on the Distri-
 457 bution and Functioning of Reservoir Microbial Communities as Determined by a Relative Distance Approach, Ecosystems, 14, 1–14,
 458 <https://doi.org/10.1007/s10021-010-9388-4>, 2011.
- 459 Weyhenmeyer, G. A. and Karlsson, J.: Nonlinear response of dissolved organic carbon concentrations in boreal lakes to increasing tempera-
 460 tures, Limnology and Oceanography, 54, 2513–2519, https://doi.org/10.4319/lo.2009.54.6_part_2.2513, 2009.
- 461 Xenopoulos, M. A., Barnes, R. T., Boodoo, K. S., Butman, D., Catalán, N., D’Amario, S. C., Fasching, C., Kothawala, D. N., Pisani,
 462 O., Solomon, C. T., Spencer, R. G. M., Williams, C. J., and Wilson, H. F.: How humans alter dissolved organic matter composition
 463 in freshwater: Relevance for the Earth’s biogeochemistry, Biogeochemistry, 154, 323–348, <https://doi.org/10.1007/s10533-021-00753-3>,
 464 2021.
- 465 Zhang, D., Shi, K., Wang, W., Wang, X., Zhang, Y., Qin, B., Zhu, M., Dong, B., and Zhang, Y.: An optical mechanism-based
 466 deep learning approach for deriving water trophic state of China’s lakes from Landsat images, Water Research, 252, 121 181,
 467 <https://doi.org/10.1016/j.watres.2024.121181>, 2024.
- 468 Zhang, Y., Yao, X., Wu, Q., Huang, Y., Zhou, Z., Yang, J., and Liu, X.: Turbidity prediction of lake-type raw water using random
 469 forest model based on meteorological data: A case study of Tai lake, China, Journal of Environmental Management, 290, 112 657,
 470 <https://doi.org/10.1016/j.jenvman.2021.112657>, 2021.