Causal Drivers of Land-Atmosphere Carbon Fluxes from Machine Learning Models and Data

Mozhgan A Farahani¹ and Allison Eva Goodwell¹

¹University of Colorado Denver

September 13, 2023

Abstract

Interactions among atmospheric, root-soil, and vegetation processes drive carbon dioxide fluxes (Fc) from land to atmosphere. Eddy covariance measurements are commonly used to measure Fc at sub-daily timescales and validate process-based and datadriven models. However, these validations do not reveal process interactions, thresholds, and key differences in how models replicate them. We use information theory-based measures to explore multivariate information flow pathways from forcing data to observed and modeled hourly Fc, using flux tower datasets in the Midwestern U.S. in intensively managed corn-soybean landscapes. We compare Multiple Linear Regressions (MLR), Long-Short Term Memory (LSTM), and Random Forests (RF) to evaluate how different model structures use information from combinations of sources to predict Fc. We extend a framework for model predictive performance and functional performance, which examines the full suite of dependencies from all forcing variables to the observed or modeled target. Of the three model types, RF exhibited the highest functional and predictive performance. Regionally trained models demonstrate lower predictive but higher functional performance compared to sitespecific models, suggesting superior reproduction of observed relationships. This study shows that some metrics of predictive performance encapsulate functional behaviors better than others, highlighting the need for multiple metrics of both types. This study improves our understanding of carbon fluxes in an intensively managed landscape, and more generally provides insight into how model structures and forcing variables translate to interactions that are well versus poorly captured in models.

Causal Drivers of Land-Atmosphere Carbon Fluxes from Machine Learning Models and Data

Mozhgan A. Farahani¹, Allison E. Goodwell^{1,2}

 $^1 \rm University$ of Colorado Denver, Department of Civil Engineering $^2 \rm Prairie$ Research Institute, University of Illinois at Urbana-Champaign

²Prairie Research Institute, 615 E. Peabody Dr. MC-650, Champaign, IL 61820

Key Points:

1

2

3

4 5

6

7

8	•	Information theory measures describe individual and joint causal relationships in
9		observed versus modeled vertical carbon dioxide fluxes.
10	•	Three machine learning models overestimate unique information from sources at
11		the expense of synergistic, or pairwise information.
12	•	Regionally trained models have improved functional performance that is not al-
13		ways captured by traditional predictive performance metrics.

 $Corresponding \ author: \ Allison \ Goodwell, \ \verb"goodwell20illinois.edu"$

14 Abstract

Interactions among atmospheric, root-soil, and vegetation processes drive carbon 15 dioxide fluxes (Fc) from land to atmosphere. Eddy covariance measurements are com-16 monly used to measure Fc at sub-daily timescales and validate process-based and data-17 driven models. However, these validations do not reveal process interactions, thresholds, 18 and key differences in how models replicate them. We use information theory-based mea-19 sures to explore multivariate information flow pathways from forcing data to observed 20 and modeled hourly Fc, using flux tower datasets in the Midwestern U.S. in intensively 21 22 managed corn-soybean landscapes. We compare Multiple Linear Regressions (MLR), Long-Short Term Memory (LSTM), and Random Forests (RF) to evaluate how different model 23 structures use information from combinations of sources to predict Fc. We extend a frame-24 work for model predictive performance and functional performance, which examines the 25 full suite of dependencies from all forcing variables to the observed or modeled target. 26 Of the three model types, RF exhibited the highest functional and predictive performance. 27 Regionally trained models demonstrate lower predictive but higher functional performance 28 compared to site-specific models, suggesting superior reproduction of observed relation-29 ships. This study shows that some metrics of predictive performance encapsulate func-30 tional behaviors better than others, highlighting the need for multiple metrics of both 31 types. This study improves our understanding of carbon fluxes in an intensively man-32 aged landscape, and more generally provides insight into how model structures and forc-33 ing variables translate to interactions that are well versus poorly captured in models. 34

³⁵ Plain Language Summary

In an agricultural landscape, exchanges of carbon dioxide between the land and at-36 mosphere occur due to photosynthesis and respiration, and depend on weather, soil, and 37 vegetation conditions. In modeling, predictive performance focuses on the relationship 38 between observed and modeled outputs, while functional performance considers the re-39 lationships between interacting inputs and outputs. We compare several performance 40 measures for three different machine learning models that simulate sub-daily carbon fluxes. 41 We look at how drivers such as solar radiation, soil moisture, temperature, humidity, and 42 rainfall provide information to carbon fluxes, and whether different machine learning mod-43 els also capture these interactions. In other words: 44

- ⁴⁵ Air, soil, and plants drive carbon's upward path,
- ⁴⁶ Models are detectives, interpreting their math.
- 47 With information theory, we map data's travel courses,
- 48 To see how models find or miss carbon's causal sources.

49 **1** Introduction

The ecohydrologic system constitutes a complex web of interactions between wa-50 ter, soil, and vegetation. The exchange of carbon dioxide (CO_2) between the land and 51 atmosphere plays a significant role in the Earth's surface temperature balance, and is 52 one of these key process affected by hydrological and ecological feedback (Liang et al., 53 2020). In terrestrial ecosystems, the carbon exchange rate is mainly controlled by the 54 photosynthesis - respiration process. Complex and nonlinear drivers such as meteorol-55 ogy, soils, vegetation, and available energy cause vertical carbon fluxes to be highly vari-56 able in space and time and challenging to measure and model (Huang et al., 2017; He 57 et al., 2018; Chen et al., 2020; Dou & Yang, 2018). Several approaches have been devel-58 oped to understand current and future terrestrial carbon flux over the past several decades 59 involving field observations (Falge et al., 2002; Xiao et al., 2011), large-scale remote sens-60 ing (Xiao et al., 2019), process-based modeling (D. Wang et al., 2011; Dunkl et al., 2021), 61 or a combination of these methods (Vetter et al., 2008; Jung et al., 2011). We take a data-62

driven approach to explore the predictability of the net CO_2 exchange rate, also known as Net Ecosystem CO_2 exchange (NEE), in agricultural landscapes in the Midwest U.S. NEE is the net carbon balance between photosynthetic CO_2 gain and respiratory CO_2 losses from plants and animals, and we use Fc as the nomenclature for NEE measured at an eddy covariance flux tower.

In this system, causal interactions need to be detected to understand interrelated 68 processes at multiple spatial and temporal scales (Runge et al., 2019; Bollt et al., 2018) 69 From a modeling perspective, this involves "intervening" in the system and manipulat-70 71 ing model structures, parameters, or inputs, and observing the resulting model behavior relative to observations (Goodwell et al., 2020). Specifically, a causal model evalu-72 ation framework should consider dependencies between inputs or source variables and 73 the target, or the "functional performance" relative to observed interactions (Goodwell 74 & Bassiouni, 2022; Bassiouni & Vico, 2021; Ruddell et al., 2019). This is particularly cru-75 cial for machine learning and deep learning models, where relationships between inputs 76 and outputs are not transparent. Understanding how these models learn, or fail to learn, 77 the dependencies we observe in nature to predict an output is vital (Goodfellow et al., 78 2016). Meanwhile, predictive performance measures capture features of the relationship 79 between the observed and modeled target output variable. In this study, we focus on the 80 functional and predictive performance of data-driven models of hourly Fc. 81

Information theory (IT) measures, which characterize uncertainty and reductions 82 in uncertainty based on probability distributions (Cover & Thomas, 2012; Shannon, 1948), 83 have been employed in various geoscience contexts to measure complexity, dependencies, 84 and driving or causal mechanisms (Balasis et al., 2013). Previous applications charac-85 terized ecohydrological process networks that reveal ecosystem behaviors (Ruddell & Ku-86 mar, 2009a; Franzen et al., 2020; Goodwell & Kumar, 2017; Ruddell et al., 2019; Sendrowski 87 & Passalacqua, 2017). Recent applications of IT-based measures in hypothesis testing 88 frameworks (Nearing et al., 2016, 2018) and to evaluate the functional performance of 89 models based on a selection of sources (Sendrowski et al., 2018; Ruddell et al., 2019; Ten-90 nant et al., 2020; Moges et al., 2022; Bassiouni & Vico, 2021; Goodwell & Bassiouni, 2022) 91 have shown great potential to better understand how models capture causal interactions 92 in various Earth systems. However, these studies tend to consider a small subset of sources 93 or a single modeled process. In this study, we take a more comprehensive view of com-94 plex ecohydrologic models and analyze information flow through the entire model. This 95 allows for identification of potential sources of model error and insights into the relation-96 ships between different components of the model. This can lead to a better understand-97 ing of the model's behavior and performance, and ultimately, more accurate predictions 98 of ecological and hydrological processes. 99

ML techniques have shown to be more effective and adaptable relative to mecha-100 nistic or semi-empirical model approaches, providing a complementary strategy to pre-101 dict carbon fluxes at local to global scales (Dou & Yang, 2018; Dou et al., 2018). Ma-102 chine learning (ML) algorithms construct empirical models based on the patterns con-103 tained in data and are very data adaptive because no assumption and functional forms 104 need to be prescribed (Jung et al., 2011). ML has been used for interpolation for gap-105 filling carbon flux data and climatic driving factors based on flux tower measurements 106 (Moffat et al., 2007; Ooba et al., 2006), decreasing the predictive errors of carbon fluxes 107 from the land surface models (T. Wang et al., 2012), and upscaling carbon fluxes of ter-108 restrial ecosystems from site to regional and global scales (Papale et al., 2015). Several 109 studies similarly indicate the ability of ML to reproduce complex ecohydrological pat-110 terns, particularly in relation to flux tower measurements (Q. Zhou et al., 2019; Tramon-111 tana et al., 2020; Reichstein et al., 2019). Specifically, Q. Zhou et al. applied a ML ap-112 proach to estimate NEE using variables such as the fraction of photosynthetically active 113 radiation (PAR), leaf area index (LAI), soil moisture, downward solar radiation, precip-114 itation, and mean air temperature. Tramontana et al. developed an ANN model to es-115

timate NEE based on the light-use efficiency concept and used a comprehensive dataset
 of soil and micrometeorological variables as flux drivers.

While machine learning models tend to make better predictions than traditional 118 models, they are often not trusted by the hydrologic community due to their black-box 119 nature (Welchowski et al., 2022). By characterizing information flow pathways and com-120 paring models beyond predictive performance, we can gain insights into their process rep-121 resentations (Goodwell & Bassiouni, 2022). This is particularly important when using 122 a certain model to extrapolate in an unknown future climate, where a model with bet-123 ter process representations may be more trustworthy to apply to an unseen scenario. In 124 this paper, we apply our IT-based model evaluation framework to three ML models, Long 125 Short Term Memory (LSTM), Random Forest (RF), and multiple linear regression (MLR) 126 to characterize how these models reproduce observed dependencies in terms of individ-127 ual, pairwise and more multivariate interactions to predict sub-daily Fc. Recurrent Neu-128 ral Networks (RNN) with LSTM are deep learning models that can successfully learn 129 long-range temporal dependencies between time steps of sequence data (Hochreiter & 130 Schmidhuber, 1997a; Sutskever et al., 2014; Kratzert et al., 2018, 2019). Meanwhile, the 131 RF is a classical ML method that is known for its capacity to handle large datasets, re-132 sist the negative impacts of noise and overfitting (Breiman, 2001), and rank the signif-133 icance of input variables (Leroux et al., 2017; Meng et al., 2021). RFs have been exten-134 sively applied in ecological classification and regression tasks (Meyer et al., 2019; Reitz 135 et al., 2021; Q. Zhou et al., 2019). We use MLR as a simple model with which to com-136 pare the more complex ML models. We develop both locally and regionally trained mod-137 els to compare model responses to larger training datasets that span multiple sites. 138

This paper is organized as follows. Section 2 describes the study site, datasets used, machine learning model development, and model evaluation. Section 3 presents the results of MLR, RF, and LSTM models. Section 4 provides a discussion, and Section 5 is a conclusion.

¹⁴³ 2 Materials and Methods

144

2.1 Site Description and Data

The data for this study was collected from multiple flux tower sites in maize/soybean 145 landscapes in the Upper Midwest Corn Belt. The Goose Creek flux tower in central Illi-146 nois (Figure 1a) is part of the NSF-funded Critical Interface Network (CINet) project 147 (https://cinet.ncsa.illinois.edu/), and collects 15-minute fluxes and meteorolog-148 ical variables at a 25m height, along with vegetation and soil properties. The Goose Creek 149 site has been extensively studied using Lidar topography and high-resolution modeling 150 of nutrient and carbon fluxes (Yan et al., 2019; Dutta et al., 2017; Woo & Kumar, 2017), 151 and footprint modeling has been applied to study how landscape heterogeneity influences 152 evapotranspiration fluxes (Hernandez Rodriguez et al., 2023). For this study, the 15-minute 153 data was resampled to hourly resolution to match with other sites. 154

We also use data from 5 maize-soybean rotation sites in the FLUXNET2015 (Pastorello 155 et al., 2020) dataset (Table 1), which provides over 1500 site-years of quality-controlled 156 datasets for various landscapes. We used the AmeriFlux version of the hourly carbon 157 flux data and meteorological variables for sites US-Ne1 (Mead - irrigated continuous maize 158 site), US-Ne2 (Mead - irrigated maize-soybean rotation site), and US-Ne3 (Mead - rain-159 fed maize-soybean rotation site). These sites are located within 1.6 km of each other at 160 the University of Nebraska Agricultural Research and Development Center near Mead, 161 Nebraska. Additionally, we used the hourly measurements of sites US-Br1 and US-Br3, 162 located in adjacent maize and soybean fields in central Iowa. The farming systems, as-163 sociated tillage, and nutrient management practices for maize/soybean production at these 164 sites are typical of those throughout the Upper Midwest Corn Belt. 165



Figure 1: (a) At a 25m height eddy covariance flux tower in Central Illinois, observed fluxes originate from up to a 10km surrounding region, dominated by a patchwork of maize and soybean fields. (b) Three flux tower sites are located in maize/soybean systems.

Table 1: Characteristics of flux tower sites. MAT, (°C) is Mean Annual Temperature. MAP (mm) is Mean Annual Precipitation.

Site ID	Name	MAT	MAP	Year	Reference
US-Ne1	Mead-irrigated contin- uous maize	10.07	790.37	2010-2021	(Suyker, 2022a)
US-Ne2	Mead-irrigated maize- soybean rotation	10.08	788.89	2010-2021	(Suyker, 2022b)
US-Ne3	Mead-rainfed maize- soybean rotation	10.11	783.68	2010-2021	(Suyker, 2022c)
US-Br1	Brooks Field Site 10-Ames	8.95	842.33	2005-2011	(Prueger & Parkin, 2016a)
US-Br3	Brooks Field Site 11-Ames	8.9	846.6	2005-2011	(Prueger & Parkin, 2016b)
CINet-GC	Goose Creek flux tower	10	900	2016-2020	(Hernandez Rodriguez et al., 2023)

The forcing variables selected for this study (Table 2) are expected to influence the dynamics of Fc between the land and atmosphere, through direct or indirect influence on photosynthesis, respiration, and other biogeochemical processes. Specifically:

169

170

171

172

173

174

175

• Ta and TS: Soil and air temperatures influence both photosynthetic rates and microbial respiration. For example, it has been found that plant respiration increases more than photosynthesis as temperature rises, which indicates that a substantial temperature increase could turn an ecosystem from a carbon source to a sink (X. Zhou et al., 2012). Meanwhile, other studies have determined that this relationship is more complex when aspects such as changing rainfall and atmospheric CO_2 concentrations are considered (Drewry et al., 2010a, 2010b; Le et al., 2011).



Figure 2: Diurnal cycle (left panel) and diurnal standard deviation cycle (right panel) of air temperature (Ta), photosynthetic photon flux density (PPFD), soil water content (SWC)) and carbon flux (Fc) over the study years corresponded to different sites (Ne1, Ne2, Ne3, Br1, Br3, GC). Each site is represented by a unique color.

Variable Description	Symbol	Unit
Carbon dioxide (CO_2) flux	Fc	$\mu molCO_2/m^2s$
Relative humidity	RH	%
Air temperature	Ta	$^{\circ}C$
Wind speed	WS	m/s
Atmospheric pressure	Pa	kPa
Precipitation	P	mm
Net radiation	NETRAD	W/m^2
Incoming photosynthetic photon flux density	PPFD *	$\mu molPhotons/m^2s$
Soil water content (volumetric)	SWC	%
Soil temperature	TS	$^{\circ}C$

Table 2: The full suite of variables used in this study.

* PAR: Photosynthetically Active Radiation $(\mu mol/m^2 s)$ in the CINet-GC site

• *RH*: Humidity levels can impact plant transpiration and stomatal conductance, 176 thereby influencing carbon uptake during photosynthesis. 177 • P and SWC: Water availability affects photosynthesis, and scarcity can lead to 178 stress conditions, slowing down carbon sequestration. 179 - PPFD and NETRAD: These radiation variables influence the energy balance and 180 are related to the amount of light available for photosynthesis, which is a primary 181 driver for carbon uptake in plants. 182 • WS: While not a direct factor, wind speed can affect plant transpiration rates, hu-183 midity levels, and even the mixing of carbon dioxide in the atmospheric layer. 184 • *Pa*: Changes in atmospheric pressure can impact gas exchange rates, indirectly 185 affecting Fc. 186

We undertook rigorous data pre-processing (SI section S1) to ensure the reliability of our analysis. This involved applying quality control measures to all datasets, and identifying and removing any outliers or erroneous patterns. We encountered missing values in some datasets, which we imputed using time series imputation methods. We note that imputation is based on certain assumptions and can introduce uncertainty, which is discussed along with the results.

193

2.2 Model Development and Experimental Design

In this study, we develop three ML models to predict Fc: Multiple Linear Regression (MLR), Long Short Term Memory (LSTM), and Random Forest (RF). Each of these models offers unique advantages and capabilities. To ensure efficient learning, all input driving variables and the output (Fc) data were normalized by subtracting the mean and dividing by the standard deviation (Minns & Hall, 1996). The output of all ML models was retransformed using the normalization parameters to obtain the final Fc prediction.

The setup of ML models necessitates the optimization of hyperparameters, a task 201 we performed via a combination of grid search and cross-validation techniques. Grid search 202 encompasses defining a range of possible parameter values and evaluating the model's 203 performance for each combination. Cross-validation helps to evaluate the model's generalization ability by partitioning the data into training and validation sets. We used a 205 5-fold cross-validation approach to search over the hyperparameter grid, where the data 206 were split into 5 subsets of equal size, and each subset was used once for validation while 207 the remaining 4 subsets were used for training. This process was repeated multiple times 208 with different partitions to ensure a robust estimate of the model's performance. 209

The ML architectures (refer to SI, Table S1) used in this study worked well for all sites in comparison to observation and were therefore chosen to be applied here without further tuning. However, a systematic sensitivity analysis of the effects of different hyperparameters was not performed in our study and could be explored in more detail in terms of their effect on predictive and functional performance.

215

2.2.1 Multiple Linear Regression Model

MLR assumes a linear function of the independent variables to predict the dependent variable. The simplicity, interpretability, and ease of use of MLR make it a popular choice for many applications. However, it assumes a linear relationship between the dependent and independent variables and is sensitive to outliers and multicollinearity. In our study, MLR provides a baseline for comparison with the more complex RF and LSTM models. We adopted the Ordinary Least Squares (OLS) method for model fitting, which optimizes the model by minimizing the sum of the squared residuals.

223

2.2.2 Random Forest Model

The Random Forest (RF) model is a powerful ensemble learning algorithm that gen-224 erates predictions by combining the outputs of multiple decision trees. Each of these trees 225 is constructed using a randomly selected subset of the features and data samples, which 226 helps to prevent overfitting. The final prediction is then derived by averaging the out-227 puts from all the trees. In a decision tree, each node represents a feature in our data, each 228 branch represents a decision rule, and each leaf represents an outcome. The root node, 229 the topmost node in a tree, corresponds to the best predictor. Decisions are made by walk-230 ing down the tree from the root to a leaf node. 231

The RF model is highly regarded for its accuracy, resilience to noise and outliers, 232 and its ability to handle high-dimensional data with nonlinear relationships and miss-233 ing values (Breiman, 2001), making it a suitable choice for our study to predict Fc. How-234 ever, due to its complexity, interpreting the model can be challenging, and the compu-235 tational cost can increase significantly with the number of trees in the forest. The per-236 formance of the RF model is significantly influenced by the fine-tuning of hyperparam-237 eters. The n-estimators (set to 100 in this study) parameter represents the number of 238 trees in the forest and a trade-off between computation time and model performance. The 239 max-depth parameter (set to 9, total number of features) controls the complexity of the 240 model, playing a crucial role in preventing overfitting. The max-features parameter (set 241 to 3), denoting the number of features to consider at each split (the maximum depth of 242 each tree), can significantly impact the model's performance and is typically set to the 243 square root of the total number of features. It is also worth noting that the random-state 244 (set to 42) parameter ensures the consistency and reproducibility of our results. 245

246

2.2.3 Long Short Term Memory Model

LSTM is a specialized form of the Artificial Recurrent Neural Network (RNN) architecture, which is designed to remember long-term dependencies in sequential data. This

capability is achieved through a unique arrangement of memory cells and three types of 249 gates: the input gate, output gate, and forget gate. These components work together to 250 selectively retain or discard information over time, making LSTM particularly adept at 251 time-series prediction tasks (Hochreiter & Schmidhuber, 1997b). We choose LSTM for 252 its capacity to model temporal dependencies in time series data, a vital characteristic 253 for accurate carbon flux prediction. We operate the LSTM in sequence-to-sequence mode, 254 in which any length of input sequence generates an equally long output sequence. We 255 chose a constant sequence length of 12 hourly time steps. This is based on the diurnal 256 cycle of environmental patterns, including temperature and light, that significantly af-257 fect Fc (Figure 2). 258

The design and training of LSTM models necessitate careful selection of various 259 parameters. These include the number of layers in the network, the number of hidden 260 units per layer, the learning rate, and the sensitivity of back-propagation to residuals be-261 tween predicted and observed outputs. Additionally, the presence or absence of dropout 262 layers, which help prevent overfitting, must be considered. To find an optimal model ar-263 chitecture, we conducted a series of experiments at different sites, manually adjusting different architectures (e.g., one or two LSTM layers or 5, 10, 15, or 20 cell/hidden units). 265 The chosen architecture consists of a two-layer LSTM network, with each layer having 266 a cell/hidden state length of 9, as number of driving source variables (Table 2). Dropout 267 layers are added between the LSTM layers to prevent overfitting (Srivastava et al., 2014), 268 and a regression layer with a single unit is added for the target variable (Fc). 269

During the training of LSTMs, each iteration step typically works with a subset 270 (called a batch or mini-batch) of the available training data. In our case, the batch size 271 is defined to be 128, and each sample in the batch consists of the Fc value and the driv-272 ing variables of the 12 preceding time steps. The loss function, calculated as the aver-273 age of the Mean Squared Error (MSE) of simulated and observed Fc of these 128 sam-274 ples, is computed in every iteration step. For faster convergence, it is advantageous to 275 have random samples in one batch. In traditional ecohydrological model calibration, the 276 number of iteration steps defines the total number of model runs performed during cal-277 ibration. The corresponding term for neural networks is called an "epoch", which is de-278 fined as the period in which each training sample is used once for updating the model 279 parameters. For instance, if the dataset consists of 1000 training samples and the batch 280 size is 10, one epoch would consist of 100 iteration steps. 281

282 2.

2.2.4 Experimental Setup

Our experimental design involves two main experiments aimed at evaluating the performance of our ML models in predicting Fc.

Local models for each site: This experiment tests the general ability of our MLMs to predict Fc at individual sites. We trained separate models for each site (Table 1) using the first 80% of the studied years as training data and the last 20% of studied years as the testing period. This resulted in six separately trained networks, one for each site.

Regional model: We train a regional model on a large dataset with data from all sites, 289 to learn general patterns and relationships between input and output data. In this, we 290 grouped all sites for the definition of the study region and used the combined data of 80%291 randomly selected for the entire period of all sites. We then test the model on each of 292 the sites separately. The regional experiment is motivated by the idea that deep learn-293 ing models perform better when trained with large amounts of data (Hestness et al., 2017; 294 295 Schmidhuber, 2015) and regional models could be a potential solution for prediction in sites without flux tower measurements (Hrachowitz et al., 2013; Sivapalan, 2003). Hav-296 ing a large training dataset allows the model to learn more generalized and abstract pat-297 terns and relationships between input and output data. For instance, if two sites behave 298 similarly, but one lacks high precipitation events or extended drought periods in the cal-299

ibration period, while having these events in the validation period, the ML model can
 learn the response behavior to those extremes and use this knowledge in the first site.

302 2.3 Model Evaluation Framework

We gauge model performance both in terms of predictive accuracy and ability to 303 encapsulate functional relationships. In this context, we consider two types of performance 304 measures: predictive performance, which assesses the model's ability to accurately pre-305 dict outcomes, and functional performance, which evaluates the model's ability to cap-306 ture the underlying functional relationships between variables (Nearing et al., 2020; Good-307 well & Bassiouni, 2022; Bassiouni & Vico, 2021). Predictive performance metrics include 308 quantitative measures of the discrepancy between the model's predictions and the actual 309 values, while functional performance can be assessed using various methods, including 310 sensitivity analysis, partial dependence plots, and information-theoretic measures. We 311 use a combination of several predictive and functional performance measures to evalu-312 ate the performance of ML models at different granularities. 313

314 2.3.1 Predictive Performance

We use Nash-Sutcliffe Efficiency (NSE) (Nash & Sutcliffe, 1970), Kling-Gupta Efficiency (KGE) (Gupta et al., 2009), and Shannon Entropy (H) (Shannon, 1948), an information theory (IT)-based measure to evaluate model predictive performance. Both *NSE* and *KGE* are widely recognized in hydrology for their effectiveness in assessing the quality of modeled predictions in relation to observed data. On the other hand, the entropy metric quantifies the uncertainty inherent in the model's predictions relative to observations. These metrics provide different perspectives on prediction errors.

The NSE is a normalized statistic that quantifies the relative magnitude of the residual variance, often referred to as "noise", in comparison to the variance of the measured data, or "information" (Nash & Sutcliffe, 1970). It is computed as follows:

$$NSE(y, \hat{y}) = 1 - \frac{\sum_{i=1}^{n} (y_i - \hat{y}_i)^2}{\sum_{i=1}^{n} (y_i - \overline{y})^2}$$
(1)

where *n* is the number of observations, $\overline{\hat{y}}$ is the mean of modeled values and y_i and \hat{y}_i are the observed and modeled values, respectively. The NSE ranges from $-\infty$ to 1. An *NSE* of 1 signifies a perfect match between modeled and observed data. An *NSE* of 0 indicates that the model's predictions are as accurate as the mean of the observed data. A negative *NSE* occurs when the observed mean is a better predictor than the model.

 $_{330}$ The *KGE* is defined by the following equation:

$$\mathrm{KGE}(y,\hat{y}) = 1 - \sqrt{(r(y,\hat{y}) - 1)^2 + (\alpha(y,\hat{y}) - 1)^2 + (\beta(y,\hat{y}) - 1)^2},$$
(2)

where r is the Pearson correlation coefficient between the observed (y_i) and modeled values (\hat{y}_i) , defined as:

$$r(y,\hat{y}) = \frac{\sum_{i=1}^{n} (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sqrt{\sum i = 1^n (y_i - \bar{y})^2 (\hat{y}_i - \bar{\hat{y}})^2}}$$
(3)

Here, n is the number of observations, and \overline{y} and $\overline{\hat{y}}$ are the mean of observed and modeled values, respectively. The variability ratio, α , is the ratio of the standard deviation

of modeled values $(\sigma_{\hat{y}})$ to observed values (σ_y) . β , the bias ratio, is the ratio of the mean

of modeled values (\hat{y}) to observed values (\bar{y}) . Similar to NSE, *KGE* values range between - ∞ and 1, where 1 represents a perfect fit.

The NSE and KGE can be more or less suitable depending on the characteristics 338 of the data and the objectives of the model (Knoben et al., 2019). NSE is based on the 339 mean squared error and is particularly sensitive to the ability of the model to reproduce 340 the variance of the data around its mean. Consequently, a model's consistent over- or 341 underestimation can influence the NSE value. If the model consistently over- or under-342 estimates the data, this will strongly affect the NSE. On the other hand, KGE also in-343 cludes the correlation between observed and simulated data in addition to bias and vari-344 ability. This enables KGE to adeptly identify patterns of over- or underestimation. More-345 over, the breakdown of the KGE into its components can provide valuable insights into 346 the model's strengths and weaknesses. A model might have a high KGE, but a low NSE347 if it reproduces the overall dynamics of the data (which KGE assesses) well but fails to 348 capture the variance around the mean (which NSE emphasizes) accurately. Conversely, 349 a model might have a high NSE, indicating a good reproduction of the observed data's 350 variance, but a low *KGE* if there are biases or variability issues. 351

IT is based on Shannon Entropy (Shannon, 1948), $H(X) = -\sum p(x) \log_2 p(x)$, where p(x) is a probability distribution function (pdf). H(X) is a measure of uncertainty of the random variable X, or the missing information that would lead to its full predictability. Here we consider the normalized difference in entropy between observed and modeled Fc as another predictive performance measure:

$$A_H = 1 - \frac{H(Fc_{mod})}{H(Fc_{obs})} \tag{4}$$

 A_H indicates how well the model captures the uncertainty that exists in the observed 357 Fc and it ranges from $-\infty$ to 1. The values of $A_H = -\infty$ never occurs in this case 358 as $H(Fc_{obs}) \neq 0$. $A_H = 0$ represents the "best" performance where the model ex-359 actly replicates the observed uncertainty. Positive values of A_H indicate that the mod-360 eled entropy $(H(Fc_{mod}))$ is lower than the observed entropy $(H(Fc_{obs}))$. In other words, 361 the model output is less uncertain, or more predictable, than the observed data. Con-362 versely, negative values of A_H indicate that the model's outputs are more uncertain than 363 the observed data. To compute pdfs, we discretize observed and modeled variables in 364 N = 100 equally sized bins spanning the minimum and maximum values of observed 365 output data. 366

367

2.3.2 Functional Performance

We also use IT to quantify the information shared between forcing variables, model 368 outputs, and observations, which can be interpreted as a measure of the model's func-369 tional performance (Nearing et al., 2020). This perspective shifts the focus from uncer-370 tainty quantification to information quantification. We explore how various model types 371 use information from driving variables (Table 2) to predict an output, or "target" vari-372 able, which here is Fc. The functional performance of a model indicates the extent to 373 which this information use is similar to or different from observed dependencies. We take 374 a multi-level IT-based approach to evaluate the functional performance of our models. 375 We will characterize complex process linkages between forcing variables or other avail-376 able information sources and Fc to assess the model's ability in capturing the relation-377 ships between the driving variables and the target variable. We consider functional per-378 formance at several different levels, specifically for individual source-target relationships, 379 pairs of sources, and all combinations of sources, or the whole model level. 380

For an individual source (X, here a forcing variable), and target (Y, here Fc), we consider reductions in uncertainty, or gains in information, in the form of mutual information as follows:

$$I(X;Y) = \sum p(x,y) \log_2\left(\frac{p(x,y)}{p(x)p(y)}\right) = H(X) - H(X|Y)$$
(5)

where I(X;Y) measures the reduction in uncertainty Y given the knowledge of X with units of bits. I(X;Y) is symmetric with respect to X and Y, and for independent variables, I(X;Y) = 0, while for fully dependent variables, I(X;Y) = min[H(X), H(Y)]. In other words, mutual information is upper bounded by the minimum uncertainty of variables involved. We calculate functional performance for individual sources based on mutual information as follows:

$$I_n(X;Z) = \frac{I(X;Z)}{H(Z)}$$

$$A_{f,MI} = 1 - \frac{I_n(X;Fc_{mod})}{I_n(X;Fc_{obs})}$$
(6)

where $I_n(X;Y)$ is the normalized MI, H(Z) is the three entropy of the target variable (Fc), 390 $I_n(X; Fc_{obs})$ and $I_n(X; Fc_{mod})$ are normalized MI of observed and modeled target vari-391 able (Fc) respectively. This captures the extent to which modeled mutual information 392 matches that of the observed target variable. $A_{f,MI}$ value close to zero represents the 393 "best" performance where the model most closely replicates the observed mutual information. 394 This can be used to assess how a model may be overestimating (negative $A_{f,MI}$ value) 395 or underestimating (positive $A_{f,MI}$ value) the influence of certain drivers, and identify 396 the most important drivers to include in a model. 397

In a more multivariate context, transfer entropy (TE) and partial information de-398 composition (PID) have been used to characterize interactions at different scales (Goodwell 399 et al., 2020). TE (Schreiber, 2000) is a specific instance of conditional mutual information, 400 which quantifies the information transferred to a target, Y_t , from a sequence of histor-401 ical states of another variable, given the knowledge of its own past states. In hydrologic 402 modeling research, TE has been used to validate and diagnose missing process connec-403 tions in a delta model (Sendrowski et al., 2018), evaluate a multi-hypothesis ecohydro-404 logical modeling framework (Bennett et al., 2019), select time aggregations and lags to-405 ward ML applications (Tennant et al., 2020), and characterize the functional performance 406 of a multi-layer canopy model (Ruddell et al., 2019). However, a TE-based analysis only 407 highlights pairwise causal connections and does not address the feature of joint or simul-408 taneous forcing from multiple drivers. Instead, we use PID to to characterize joint in-409 fluences from multiple source variables to a target (Williams & Beer, 2010; Goodwell et 410 al., 2020). For example, previous studies have compared how stomatal optimization mod-411 els respond to soil water supply and atmospheric demand (Bassiouni & Vico, 2021), how 412 simple to complex models behave under different source dependencies (Goodwell & Bassiouni, 413 2022), and stomatal model representations of physiological limits on transpiration (Hawkins 414 et al., 2022). We consider two sources, or model forcing variables, that provide information 415 to a target variable, which could be an observation or a model output. In a system where 416 two sources share information from X and Y with a target Z, the total information quan-417 tity, I(X, Y; Z), can be partitioned into synergistic (S), unique (U), and redundant (R) 418 components. This partitioning is as follows: 419

$$I(X,Y;Z) = S_{X,Y} + R_{X,Y} + U_{X|Y} + U_{Y|X}$$
(7)

Here, $S_{X,Y}$ is synergistic information or joint information that is provided only when both 420 sources are known together. $R_{X,Y}$ is redundant information or overlapping information 421 that both sources provide individually. $U_{X|Y}$ and $U_{Y|X}$ terms indicate unique information 422 that individuals influence when one source provides information that is not provided by 423 the other. We use a partitioning method described in Goodwell and Kumar to obtain 424 these components of the total information (refer to SI section S2 for more details). We 425 normalize components by dividing each by the total mutual information I(X, Y; Z), such 426 that all information components add up to 1, and a given component indicates the frac-427 tion of reduced uncertainty in Z that can be attributed to that information type. These 428 IT-based measures R, U, and S characterize different types of causal relationships be-429 tween variables. They are particularly useful to interpret multivariate interactions, such 430 as the Fc-related processes of interest here. 431

For computing mutual information and information partitioning components, we 432 used different number of bins, based on the range of observed and modeled data (i.e., the 433 difference between the maximum and minimum values). We calculated the number of 434 bins for the model by taking the ratio of the range of the model to the range of the ob-435 servation, multiplied by the number of bins in the observations (N = 100). This method 436 effectively scales the number of bins based on the relative range of the model and observed 437 data, with the assumption that a wider range would need more bins to capture the data 438 distribution effectively. We compute statistical significance of observed or modeled IT 439 measures using a shuffled surrogates approach (Ruddell & Kumar, 2009b). Details on 440 these methods are provided in SI, Section S3. 441

We use PID to calculate the pairwise functional performance in terms of redundancy, synergy, and unique information and "overall" information partitioning for a given pair of sources. We consider the pairwise functional performance as the relative difference in an information flow measure for modeled versus observed data, separated into different components related to information partitioning measures S, R, and U, (Equation 7), respectively as $A_{f,S}, A_{f,R}$, and $A_{f,U}$ (Goodwell & Bassiouni, 2022). For example:

$$A_{f,S_{i,j}} = S(X_i, X_j; Z_{mod}) - S(X_i, X_j; Z_{obs}); \quad \text{for } i \neq j$$

$$\tag{8}$$

where X_i and X_j indicate two source variables. The same concept applies for R. For unique information, we consider the sum of the two unique components (U_X+U_Y) . A positive value indicates that the model overestimates a particular component at the expense of a different information type. The partitioning functional performance for a pair of sources is defined as the sum of the absolute values of the three pairwise measures as follows:

$$A_{f,Ipart_{i,j}} = |A_{f,S_{i,j}}| + |A_{f,R_{i,j}}| + |A_{f,U_{i,j}}|$$
(9)

This measure ranges from 0, for a model that exactly reproduces the observed information 453 components, to 2, for a model that entirely substitutes one type of information for an-454 other or a combination of other information types. For instance, if the observed system 455 shows that U = 1 (all information is unique), but a model system estimates S = 1 (that 456 all information is synergistic), this leads to $A_{f,S} = 1$, $A_{f,U} = -1$ and $A_{f,Ipart} = 2$. 457 While the individual source level identifies how the ranking of modeled variable impor-458 tance differs from observations, this pairwise level identifies how the model is interpret-459 ing information provided by combinations of sources. 460

461 At the highest "whole model" level of analysis, we calculate average overall func-462 tional performance across all individual $(A_{f,MI})$ and pairs of sources $(A_{f,Ipart})$ as fol-463 lows:

$$A_{f,MI,tot} = \frac{\sum_{i=1}^{n} (1 - |A_{f,MI_i}|)}{n},$$
(10)

464 and

$$A_{f,Ipart,tot} = \frac{\sum_{i=1}^{n-1} \sum_{j=i+1}^{n} (2 - A_{f,Ipart_{i,j}})}{(n^2 - n)},$$
(11)

where n is the number of source variables. $A_{f,MI,tot}$ ranges from $-\infty$ to 1 and $A_{f,Ipart,tot}$ 465 ranges from 0 to 2. We note that these measures are the originally defined individual and 466 pairwise performance measures subtracted from 1 or 2, in order to align higher values 467 with "best" model performance. In other words, a value of 1 (or 2 for $A_{f,Ipart,tot}$) now 468 corresponds to a perfect match of modeled values to the observed data (Table 3). This 469 level of functional performance metrics gauges the model's overall ability to replicate the 470 observed interactions. Figure 3 and Table 3 indicate the different levels of functional and 471 predictive performance analysis. 472

Metric	Range	Best Per- formance	Eq. No.	Description
NSE	$-\infty$ to 1	1	1	Nash-Sutcliffe Efficiency (predictive)
KGE	- ∞ to 1	1	2	Kling-Gupta Efficiency (predictive)
A_H	$-\infty$ to 1	0	4	Normalized difference in entropy be- tween observed and modeled (predic- tive)
A _{f,MI}	$-\infty$ to 1	0	6	MI difference for individual source (functional)
$\begin{array}{c} A_{f,S_{i,j}}, A_{f,R_{i,j}}, \\ A_{f,U_{i,j}} \end{array}$	-1 to 1	0	8	Information partitioning components difference for a pair of sources (func- tional)
$A_{f,Ipart_{i,j}}$	0 to 2	0	9	Overall information component differ- ence for a pair of sources (functional)
$A_{f,MI,tot}$	$-\infty$ to 1	1	10	Average functional performance of in- dividual source level across all driving sources
$A_{f,Ipart,tot}$	0 to 2	2	11	Average functional performance across all pairs of sources for overall information partitioning

Table 3: Summary of predictive and functional performance metrics.

473 **3 Results**

3.1 Predictive Performance

NSE and KGE values are higher for local relative to regional training across all
 ML models and sites (Figure 4a). This implies that local training allows the models to



Figure 3: Illustration of functional and predictive performance. Nodes represent driving sources and target variables, and arrows represent different levels of functional performance. Predictive performance (*NSE* and *KG* and A_H) measure agreement between observed and modeled values (Equations 1, 2, and 4). Blue, red, and green links show relationships that can be captured by functional performance metrics at different levels (Table 3).

better capture certain characteristics of each site. The regional model performance may
stem from the limitations of this study, mainly a relatively small number of sites and siteyears. A more extensive dataset encompassing multiple sites over varied temporal spans
may provide the model with a broader range of conditions and variability, enabling it
to generalize more effectively.

Meanwhile, we find that the A_H of local models is higher than that of regional models (Figure 4b). A negative A_H occurs when $H_{mod} > H_{obs}$. This means that regional models actually introduce greater variability or uncertainty in Fc relative to observations. It is important to note that a negative A_H does not indicate "inferior" performance, since values close to zero represent "best" performance where the models reproduce the observed H(Fc). While regional models over-estimate uncertainty in Fc, locally trained models underestimate uncertainty to a similar degree (Figure 4b).

⁴⁸⁹ When comparing performances of the three different models, RF (square markers ⁴⁹⁰ in Figure 4a) consistently exhibits higher *NSE* and *KGE* values across all sites and both ⁴⁹¹ training experiences. This indicates the robustness of the RF model irrespective of the ⁴⁹² scale of the training data. Moreover, RF generally performs well in capturing the uncer-⁴⁹³ tainty in the observed Fc in both local and regional scales (square markers, Figure 4b). ⁴⁹⁴ RF models have the best A_H performance for both regional and local models, indicat-⁴⁹⁵ ing their ability to replicate the observed entropy of Fc.

⁴⁹⁶ MLR (circle markers in Figure 4) performance varies highly between sites. For some ⁴⁹⁷ sites, the *NSE* values are very low, especially for regional training, suggesting MLR does ⁴⁹⁸ not capture the specific behaviors of those sites effectively. The negative *NSE* values in-⁴⁹⁹ dicate that a mean predictor would have been better for most sites. Meanwhile, *KGE* ⁵⁰⁰ values fall closer to the 1:1 line of Figure 4a, indicating that the *KGE* metric does not ⁵⁰¹ distinguish as many differences between regional and local training. Similarly, A_H for



Figure 4: Predictive performance, (a) NSE (filled markers) and KGE (empty markers), and (b) the normalized difference in entropy between observed and modeled values (A_H) of three different models (MLR, RF, and LSTM, marker shapes) trained on local and regional data for six different sites (Table 1). Colors denote sites. The 1:1 line indicates equal performance for local and regional models.

the MLR model has the most spread between the study sites. For Nebraska sites (Ne1, Ne2, and Ne3), MLR has negative A_H values, which suggests that MLR model's outputs for these sites are more uncertain compared to the observed data. On the other hand, MLR for the other sites show positive A_H values.

LSTM (triangle markers in Figure 4a) results in NSE and KGE values between those 506 of RF and MLR. For some sites, performance is close to that of the RF. This suggests 507 that LSTMs can model temporal patterns at individual sites to some extent, and is al-508 ways better than a mean predictor, but it never outperforms the RF model given the same 509 training data. Given that LSTMs can model temporal sequences, the varied performance 510 suggests that while some regional patterns are temporal, others might be non-sequential. 511 We also find similar behaviour for LSTM as RF in capturing the entropy of observed Fc, 512 except for more variability between sites. When models are trained locally, LSTM mod-513 els tend to produce outputs that are less uncertain, or more predictable, than the observed 514 data $(A_H > 0)$. When models are trained regionally, LSTM outputs are more uncer-515 tain than observations. This difference between local and regional training for both LSTM 516 and RF indicates that the regional training enables the model to produce more variable 517 outputs, while local training leads to a more restricted range of Fc. 518

519

3.2 Functional Performance

At the individual and pairwise level, we focus on a single site, Ne1, as the site with the highest predictive performance and few gaps in forcing variables (*WS* and *NETRAD*). Other sites show similar patterns in mutual information and information decomposition measures, and we present full results for these in the Supplementary Information (SI Figures S3-S18).

525

3.2.1 Individual Source Level

Each variable is ranked based on the average observed MI across all sites (Figure 527 5a, black line). TS and Ta share the most information with Fc, indicating a strong de-528 pendence on fluctuations in both air and soil temperatures. The next variables that share 529 information with Fc are radiation variables, NETRAD and PPFD. Meanwhile, precipitation (P) is a very weak predictor of Fc, which is expected since sub-daily precipitation contains many zero-values, leading to low entropy. Instead, we see that SWC shares more information with Fc, indicating that moisture available to roots and soil is important. Meteorological variables Pa, RH, and WS are relatively weak individual predictors. Models either overestimate or underestimate these mutual information values, resulting in a different ranking of variables for each model type (Figure 5a).

⁵³⁶ We use $A_{f,MI}$ to assess the extent to which mutual information matches with the ⁵³⁷ observed target variable at Ne1 site (Figure 5b) and at other sites (SI Figure S3-S6). Higher ⁵³⁸ absolute $A_{f,MI}$ values suggest that the modeled value is far from the observed value. If ⁵³⁹ $A_{f,MI}$ is negative, the model overestimates the mutual information of observed Fc (an ⁵⁴⁰ overly deterministic model), and if $A_{f,MI}$ is positive, the model underestimates observed ⁵⁴¹ mutual information (an overly random model).

The MLR model tends to underestimate mutual information (positive $A_{f,MI}$) for 542 TS, Ta, SWC, Pa, WS, and P while overestimating for NETRAD, and PPFD, particu-543 larly for local training (Figure 5b, blue circles). MLR also shows the largest spread in 544 over and underestimates of mutual information. The LSTM model for local training has 545 a negative $A_{f,MI}$ for the most relevant drivers, but this is improved under regional train-546 ing (Figure 5b, green triangles). The RF models closely replicate observed mutual information 547 for both regional and local training (Figure 5b, red and orange squares). This highlights 548 the power of RF in capturing the intricacies and dependencies within Fc regardless of 549 the scale of the training data. Here we discuss the model representation of individual forc-550 ing variables. 551

552	•	TS, Ta: While local and regional MLR model greatly underestimates the influence
553		of temperature variables, the locally trained LSTM model overestimates it to a
554		similar degree. In other words, the local LSTM model correctly identifies these
555		as top sources of information to Fc, but to a more extreme degree, while the MLR
556		models do not consider temperature as a top source.
557	•	NETRAD and PPFD: For local MLR, $A_{f,MI}$ is large and negative, indicating that
558		the model overestimates the influence of radiation variables, and interprets them
559		as the most important forcing variables instead of temperatures. However for the
560		regional MLR, $A_{f,MI}$ is close to zero, indicating that the regional model mitigates
561		this over-estimation. The only model that slightly underestimates mutual information
562		from these variables is the regionally trained LSTM.
563	•	SWC and P : Precipitation is a very weak driver according to both observations
564		and models (Figure 5a), but models nearly always underestimate mutual information.
565		They also underestimate information from SWC, except for the regionally trained
566		LSTM. This indicates that models may lack sensitivity to moisture variability.
567	•	WS: Across all models, the $A_{f,MI}$ values are fairly consistent, small, and positive,
568		indicating all models slightly underestimate the influence of wind speed.
569	•	Pa and RH : The locally trained MLR model shows the worst performance in terms
570		of both over and under-estimating information from these variables.

These patterns are similar for other sites and under regional training (SI Figures S3-S6). This consistency suggests that the observed MI behaviors are not merely sitespecific but possibly representative of broader environmental interactions. The key takeaway is that all models overestimate the influence of certain drivers at the expense of others, but to different degrees. This understanding can be useful to refine models or test the sensitivity of certain drivers. However, this level of analysis may omit drivers that provide information jointly rather than individually.



Figure 5: (a) Normalized mutual information (I_n) and (b) functional performance for individual variables $(A_{f,MI})$, Equation 6, for Multiple Linear Regression (MLR), Random Forest (RF), and Long Short-Term Memory (LSTM) models, under local and regional training at Ne1 site. Each variable is ranked (order on x-axis) based on the average observed MI across all sites (black line).

578 3.2.2 Pairwise and Model Level

In the observed data, most variable pairs provide synergistic (S) or unique information 579 (U) to Fc (Figure 6a-c). The only pairs that provide a large fraction of redundant information 580 (R) are closely related pairs (Ta, TS) and (PPFD, NETRAD). However, we note that 581 their redundancy is still less than 0.5 as a fraction of total information, and the other 582 half of the information they provide is U. Precipitation (P) provides the most U when 583 paired with other variables (Figure 6c), but as found in the previous analysis of individ-584 ual sources, the actual amount of information it provides is very small due to its low en-585 tropy. Meanwhile, Ta tends to provide the next highest fraction of U when paired with other sources, while RH and WS to provide S along with other sources. In general, re-587 gardless of the amount of information that sources provide, here we find that they mainly 588 provide unique and synergistic information types. 589

All models tend to underestimate S (negative $A_{f,S}$, Figure 6d,g,j) for most vari-590 able pairs, at the expense of overestimating U (positive $A_{f,U}$, Figure 6f,i,l). For exam-591 ple, in the MLR model, RH greatly underestimates S and overestimates U when paired 592 with other variables (Figure 6d, f). While the underestimation of synergistic relationships 593 is widespread, the overestimation of redundancy only tends to occur for the most cor-594 related variable pairs, specifically (Ta, TS) and (PPFD, NETRAD). This indicates that 595 models rely excessively on these correlations, which results in an overemphasis in R. In 596 other words, the observed relationship between these variables is not as redundantly in-597 formative for Fc as the model predicts, but they are instead more unique predictors. 598

Essentially, depending on the variable pair, the model either uses information uniquely 599 where observations show a synergistic type of relationship, or uses information redun-600 dantly where observations show both unique and redundant contributions. The MLR 601 model shows the largest trade-off between S and U partitioning performances (Figure 6d,f), followed by LSTM. Meanwhile, MLR is the only model that does not overestimate 603 R provided by (Ta, TS), and in fact captures all information types accurately for this 604 pair. However, we note that this MLR model also greatly underestimates the individual 605 information components shared by each of these variables to the target (Figure 5). In 606 other words, the MLR model greatly underestimates the importance of these tempera-607 ture variables as predictors of Fc, but does reflect the mechanism by which they jointly 608 provide information. 609

While broad patterns in information decomposition components are similar between models, there are several differences. For example, consider the (SWC, Ta) pair (bottom corner in all Figure 6 panels). For MLR, the information components are reproduced fairly accurately. For RF, U is overestimated at the expense of S to a minor degree. For LSTM, this occurs to a higher degree and R is also slightly overestimated. Meanwhile the MLR model greatly overestimates U from the pair (RH, NETRAD) at the expense of S, while the other two models have a similar but less extreme pattern.

When we consider the combined partitioning performance, $A_{f,Ipart}$ for each vari-617 able pair, the RF model has the best model performance, as it shows more $A_{f,Ipart}$ val-618 ues close to zero (Figure 7). The MLR shows the most variability between pairs of sources, 619 such that some pairs have very good functional partitioning performance and others have 620 values of $A_{f,Ipart}$ greater than 1, indicating that over half of the information decompo-621 sition is misrepresented by the model. RH, NETRAD, and PPFD have particularly poor 622 functional performance when combined with other sources for the MLR model. The LSTM 623 model also has lower functional partitioning performance relative to RF, but behavior 624 is more even between pairs of variables. Precipitation (P) always has the best functional 625 performance when paired with other variables, but it is the weakest source and provides 626 very little information regarding Fc for either models or observations. 627

When we consider other sites (SI Figures S7-S12), we find similar patterns in pair-628 wise functional performance, specifically the overestimation of U at the expense of S and 629 overestimation of R for correlated source pairs. However, we find that regionally trained 630 models diminish some of the issues observed in the localized models. The broader dataset 631 that regional training offers seems to provide a more balanced representation, allowing 632 models to discern patterns beyond local-specific interactions. The regional model also 633 corrects the balance between synergy and unique contributions, leading to a more accu-634 rate representation of how these variables interact. This trend is especially evident in 635 the LSTM model, which demonstrates enhanced functional performance under regional 636 training (SI Figures S13-S18). In terms of site differences, we find that regional LSTM 637 model has the best model performance at Ne1 and Ne3 sites and RF model has the best 638 performance among other sites. 639

When we calculate average overall functional performance at individual level $(A_{f,MI,tot})$, 640 we find patterns that are similar to the average pairwise functional performance $(A_{f,Ipart,tot})$ 641 (Figure 8). Specifically, local RF models perform slightly better than regional RF mod-642 els on the individual level, while regional MLR and LSTM models generally perform bet-643 ter than the local models (Figure 8a). However, at the pairwise level, regional models 644 consistently outperform their local equivalents (Figure 8b). This contrasts with trends 645 observed in the predictive performance metrics (Figure 4), where local training led to 646 higher NSE values relative to regional training. 647

Among all models, the RF model demonstrates the best performance, both at in-648 dividual and pairwise levels (square markers in Figure 8). For individual sources, local 649 RF models have better performance than the regional models. But when considering pair-650 wise relationships, the regional RF model shows superior performance. On the other hand, 651 the MLR model exhibits the lowest performance values at the individual level but per-652 forms more similarly to LSTM when considering pairwise relationships. The regional LSTM 653 model also shows good performance at both the individual and pairwise levels. However, 654 the performance of the local LSTM model varies more across different sites at the indi-655 vidual level, while the pairwise performance is more consistent for the regional model. 656 This analysis highlights that changes in one aspect of functional performance do not nec-657 essarily translate to similar changes in other aspects. 658

659

3.3 Relationship between Predictive and Functional Performance

The relationship between predictive performance and functional performance pro-660 vides insights into how a model balances replicating the observed data and its ability to 661 capture observed relationships. As an illustration, we first focus on two key metrics: the 662 KGE, representing predictive performance, and the $A_{f,Ipart,tot}$, indicating functional per-663 formance (Figure 9). For the 6 sites, two training types, and 3 model types, we have 36 664 total model runs for this comparison. All models show higher functional performance 665 under regional training, but differences in KGE are on a site-by-site basis. The Ne1 site 666 tends to be the highest performing site for all models in terms of KGE, but varies be-667 tween models for $A_{f,Ipart,tot}$. 668

The functional and predictive performances for RF are both high relative to other models, and there is little variability between sites. However, there is an apparent tradeoff between functional and predictive performance, in that sites with the highest KGE tend to have lower $A_{f,Ipart,tot}$. Meanwhile, there is a slight positive trend for locally trained LSTM and MLR models, where higher functional and predictive performances go together (Figure 9).

A correlation analysis shows that while functional and predictive performance measures tend to be correlated to each other (Figure 10a,c), there are fewer statistically significant (p < 0.05) correlations between the two types (Figure 10b). This correlation analysis is based on all 36 model cases (3 ML models, regional and local, and 6 sites) so



Figure 6: Observed pairwise (a) synergistic $(S_{i,j})$, (b) redundancy $(R_{i,j})$, and (c) uniqueness $(U_{i,j})$ information flow at Ne1 site. Pairwise functional performance of three models under local training experience at Ne1 site. The heat-map represents the relative difference in information decomposition partitioning measures $(A_{f,S_{i,j}}, A_{f,R_{i,j}}, \text{ and } A_{f,U_{i,j}})$ between modeled and observed data for each pair of forcing variables. Positive values (green) in (d)-(l) indicate that the model overestimates the information type, while negative values (red) indicate underestimations.



Figure 7: Pairwise functional partitioning performance $A_{f,Ipart_{i,j}}$ for (a) MLR, (b) RF, and (c) LSTM models under local training experience at Ne1 site. Values close to zero indicate optimal partitioning performance for a given pair.



Figure 8: The whole model functional performance for (a) across all individual sources, $A_{f,MI,tot}$ and (b) across all pairs of sources, $A_{f,Ipart,tot}$), of three model types under two training experiences, local and regional, for six sites.

does not distinguish trends for a single model type or training experience. As illustrated in Figure 9, there may be a stronger correlation within a given model type and training. We split the KGE into its three constituent components, where high values of each term indicate "best" model performance. Similarly, the A_H measure of entropy and functional performance metrics are scaled so that high values indicate best performance, and positive correlations are easy to interpret.

Predictive performance metrics are positively correlated, except for the α , or vari-685 ability, term of KGE with NSE and A_H . We find that the correlation component (r) is 686 most correlated to the total KGE. Meanwhile, β and α terms are less correlated to KGE. and individual KGE components are less correlated to each other. This indicates that 688 the correlation between observed and modeled Fc is the most predictive of KGE for these 689 models. Meanwhile, both β and r terms are highly correlated with NSE. This highlights 690 that the NSE is sensitive to the bias between model and observations and their corre-691 lation. The two full model functional performance metrics are also positively correlated 692 (Figure 10c), indicating that models with high performance in terms of individual sources 693 also reproduce pairwise relationships well. 694

In terms of correlations between functional and predictive measures (Figure 10b), 695 5 of the 12 possible correlations are positive and the other 7 are non-statistically signif-696 icant, indicating that higher predictive performance is generally but not always associ-697 ated with higher functional performance. The KGE α , or variability, component shows 698 the highest correlation with functional measures, followed by the total KGE. This leads 699 us to interpret that α is the most indicative of functional performance, and is the basis 700 for the correlation between KGE and the functional measures. This indicates that mod-701 els that reproduce the standard deviation of observed Fc, upon which α is based, also 702 tend to reproduce observed forcing-Fc relationships at both a pairwise and individual 703 level. Meanwhile, A_H , which is based on the difference in entropies of observed and mod-704 eled Fc, does not have a statistically significant correlation with functional performance. 705 This illustrates that a model can reproduce the entropy of the observation, but not re-706 produce the distribution or functional relationships. In other words, the entropy is a sum-707 mary statistic that does not necessarily indicate whether the model correctly replicates 708 other features of the distribution of the data. No functional performance measures are 709 correlated to the NSE, the β , or bias component of KGE, or A_H . This could be related 710 to the linearity of these predictive performance measures that may not reflect nonlinear 711 and joint interactions detected with mutual information. Additionally, we note that IT-712 based measures consider the distribution of the data but not the actual values, such that 713 an IT measure would not capture a constant bias between two variables. 714

715 4 Discussion

Many machine learning approaches have been applied across major sub-domains 716 of Earth system science and are increasingly being integrated into operational schemes 717 and used to discover patterns, improve our understanding, and benchmark physically-718 based models. Ideally, ML models generate predictive models devoid of any presumptions 719 on the underlying ecological structure or the mathematical representation of processes 720 and interactions in an ecosystem. However, this lack of presumptions is correlated to a 721 lack of understanding of whether and how these models are capturing functional relation-722 ships that exist in nature. The results of this study emphasize that functional performance—how 723 accurately models capture the underlying relationships between variables—can be paired 724 with more traditional metrics of model performance. By evaluating both functional and 725 predictive aspects and their interrelationship, we can obtain a wider perspective on the 726 strengths and limitations of different machine learning models. This multi-tiered approach 727 not only can be used to explore the behavioral ranges for both machine learning and process-728 based models but also guides model development by highlighting model deficiencies based 729 on information flow pathways that would not be apparent based on existing measures. 730



Figure 9: Predictive performance (KGE) and the overall model level of functional performance $(A_{f,Ipart,tot})$ of three model types under two training experiences, local (filled markers) and regional (empty markers).

Since ML-predicted fluxes can be used as benchmarks for physical land-surface and climate model evaluation (Q. Zhou et al., 2019; Anav et al., 2015; Best et al., 2015), it is
valuable to understand nuances in their behavior.

While earlier studies on the CO_2 balance of vegetated surfaces applied linear re-734 gression for estimating the carbon fluxes (Jensen et al., 1996; Xu & Qi, 2001; Burrows 735 et al., 2005), artificial neural network (ANN) and the support vector machine (SVM) meth-736 ods have also been used to estimate terrestrial carbon fluxes and interpret the nonlin-737 ear relationship between ecosystem-based carbon fluxes and environment variables based 738 on eddy covariance measurements (Papale & Valentini, 2003; Dou & Yang, 2018). For 739 example, an ANN was able to filter out noise, predict the seasonal and diurnal variation 740 of carbon fluxes, and extract patterns such as increased respiration in spring during root 741 growth, which was formerly not well represented in carbon cycle models (Papale & Valen-742 tini, 2003). In this study, the Random Forest model showed both the highest functional 743 and predictive performances, confirming that its better predictions really are associated 744 with better process representations. The RF's non-parametric nature means it makes 745 fewer assumptions about the underlying relationships between variables, thus enabling 746 it to proficiently model intricate, non-linear interactions. Meanwhile, linear regression 747 had a wide spread in performance levels between individual sites, and greatly overesti-748 mated the influence of radiation drivers that are highly linearly correlated to carbon flux. 749 The LSTM model performance varied greatly between local and regional training, indi-750 cating that its functional performance benefited from training data from multiple sites. 751

Complex and nonlinear drivers such as meteorology, soils, vegetation, and available 752 energy cause Fc to be highly variable in space and time and challenging to measure and 753 model (Huang et al., 2017; He et al., 2018; Chen et al., 2020; Dou & Yang, 2018). Sev-754 eral approaches have been developed to understand current and future terrestrial car-755 bon flux over the past several decades involving field observations (Falge et al., 2002; Xiao 756 et al., 2011), large-scale remote sensing (Xiao et al., 2019), process-based modeling (D. Wang 757 et al., 2011; Dunkl et al., 2021), or a combination of these methods (Vetter et al., 2008; 758 Jung et al., 2011). Our study sheds further light on how forcing variables provide information 759 to observed carbon fluxes. We found that temperature and radiation variables are most 760 highly informative of Fc, followed by moisture-related variables such as RH and SWC. 761 While many variables have a diurnal pattern, including Fc, we find that forcing variables 762 tend to provide synergistic or unique information, rather than redundant information, 763



Figure 10: Correlation (*p-value* < 0.05) between performance metrics listed in Table 3 (scaled so that larger values always correspond to best performance), for the 36 model runs performed in this study. (a) and (c) separate correlations within predictive and functional categories, respectively, while (b) shows correlations between functional and predictive metrics.

indicating that the overlap in information content is relatively low. Meanwhile, RH is 764 relatively weak as an individual source, but we found that it provides synergistic information 765 when paired with many other sources. This indicates that the relevance of a variable like 766 RH could be underestimated in an analysis that did not consider multivariate interac-767 tions, since it is a weak individual source but enhances the information content of other 768 sources. In terms of modeling, we found that MLR, the simplest model, overestimates 769 information from radiation variables and underestimates information from temperatures. 770 This suggests that MLR captures the strongly linear diurnal pattern between energy avail-771 ability and carbon flux, but misses a stronger but more nonlinear relationship with tem-772 perature due to the limitations in its parameterization. Finally, the tendency of all mod-773 els to underestimate information from SWC indicates that water availability to plants 774 is a complex driver of Fc that is difficult to capture in a functional form. 775

We note several limitations and assumptions that could be improved in future work. Future research could delve deeper into variations between sites, exploring what site-specific features influence model performance. One of the uncertainties of using flux tower measurements to estimate Fc is the impact of shifting land cover on the accuracy of the ob-

servations. The land-atmosphere exchange fluxes that generate carbon flux are influenced 780 by the dynamic upwind surface area, called the flux footprint, which can exhibit spatial 781 heterogeneities (Hernandez Rodriguez et al., 2023; Leclerc & Foken, 2014). As a result, 782 fluxes from different sources can mix at the observation point, introducing uncertainty 783 into the measurements. Meanwhile, this study assumes that the mix of crop types be-784 tween sites and between observation time points leads to similar causal interactions be-785 tween forcing variables and carbon flux. We also did not specifically focus on the opti-786 mization of hyperparameters within each ML model, which could have an effect on func-787 tional and predictive performance. Moreover, the precision and general quality of the 788 forcing variables and Fc are important as they have underlying uncertainties and have 789 been gap-filled, and our interpolation methods may have more effect on some model struc-790 tures than others and future research could explore how models use information encoded 791 in forcing data (Farahani et al., 2022). We also note that the MLR performance can be 792 significantly influenced by multicollinearity among the forcing variables, and we did not 793 test for this aspect. In terms of data size, we only considered six locations and approx-794 imately 50-site years, so further studies could more specifically consider the effect of in-795 creasingly large and diverse training datasets on model functional behaviors. Finally, the 796 models evaluated represent just a fraction of the available algorithms, and we do not con-797 sider a wider range of ML and process-based models. 798

While predictive and functional metrics tend to be positively correlated, there are 799 cases where a model change could be made that appears to improve predictions, but sac-800 rifices a functional relationship. For example, the finding that regionally trained mod-801 els tend to have improved functional performance indicates that these models can dis-802 cern patterns beyond local-specific interactions. However, in this study the predictive 803 performance of regional models was somewhat lower relative to single-site models, po-804 tentially marking a trade-off between functional and predictive performance. A "perfect" 805 model should replicate all functional relationships as they are observed, but it still may 806 not have perfect predictive performance due to missing information. In other words, the 807 forcing variables simply do not contain all the information necessary to make a perfect 808 accurate prediction. In this way, information-based metrics of functional performance 809 provide a type of upper bound for predictive performance. This underscores the need 810 for a nuanced approach to model selection. For an ungauged site with no validation data, 811 a regionally trained model is likely the most applicable since it has a stronger functional 812 performance and can reproduce processes as they are observed. The LSTM model was 813 the most responsive to changes in training data size, which could relate to its complex-814 ity and need for many datasets to learn time-dependent interactions. 815

⁸¹⁶ 5 Conclusion

Predictive accuracy is just one facet of modeling complex ecohydrologic systems. 817 Meanwhile, functional performance metrics capture how a model grasps the intricate re-818 lationships among variables. In order to use models for prediction in unseen conditions, 819 and compare between machine learning and physically based model structures, we need 820 to ensure that models don't just predict well, but also understand and represent the un-821 derlying processes effectively. In other words, understanding the why and how behind 822 predictions can be as vital as the predictions themselves. In this study, the Random For-823 est model emerged as a consistently reliable model in terms of both predicting carbon 824 fluxes and reproducing observed functional relationships at multiple levels. Meanwhile, 825 a simple linear regression will overestimate the influence of variables with the most lin-826 ear relationships to the target outcome. All models in this study had the common fea-827 ture of underestimating synergistic interactions and overestimating unique ones. This 828 indicates that all models are not quite capturing information flows at higher levels, where 829 multiple sources provide information to the target jointly, and indicates that even the 830 models with the highest predictive performance could be improved. Similarly, while per-831

formance measures tend to be correlated, no single performance measure captures the

effect of all the others. This study advocates for a combined approach to model evaluation and validation, which considers both predictive performance and how the model

ation and validation, which considers both predictive performance

captures interactions in the ecohydrologic system.

836 Acknowledgments

- M.A. Farahani and A.E. Goodwell acknowledge funding from NSF Grant EAR #2012850
- for the Critical Interface Network for Intensively Managed Landscapes (CINet) and the
- NASA New Investigator Grant #80NSSC21K0934. Python codes for analyses presented
- here are available on at https://github.com/allisongoodwell/Farahani_CarbonML2023.

841 **References**

- Anav, A., Friedlingstein, P., Beer, C., Ciais, P., Harper, A., Jones, C., ... Zhao, M.
 (2015). Spatiotemporal patterns of terrestrial gross primary production: A
 review. *Reviews of Geophysics*, 53(3), 785-818. doi: https://doi.org/10.1002/
 2015RG000483
- Balasis, G., Donner, R. V., Potirakis, S. M., Runge, J., Papadimitriou, C., Daglis,
 I. A., ... Kurths, J. (2013). Statistical mechanics and information-theoretic
 perspectives on complexity in the Earth system (Vol. 15) (No. 11). doi:
 10.3390/e15114844
- Bassiouni, M., & Vico, G. (2021). Parsimony versus predictive and functional performance of three stomatal optimization principles in a big-leaf framework. New Phytologist, 0–2. doi: 10.1111/nph.17392
- Bennett, A., Nijssen, B., Ou, G., Clark, M., & Nearing, G. (2019). Quantifying process connectivity with transfer entropy in hydrologic models. *Water Resources Research*, 55(6), 4613-4629. doi: 10.1029/2018WR024555
- Best, M. J., Abramowitz, G., Johnson, H. R., Pitman, A. J., Balsamo, G., Boone,
 A., ... Vuichard, N. (2015). The plumbing of land surface models: Benchmarking model performance. *Journal of Hydrometeorology*, 16(3), 1425 - 1442. doi: https://doi.org/10.1175/JHM-D-14-0158.1
- Bollt, E. M., Sun, J., & Runge, J. (2018). Introduction to focus issue: Causation inference and information flow in dynamical systems: Theory and applications. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 28(7), 075201. doi: 10.1063/1.5046848
- Breiman, L. (2001). Random forests. Machine Learning, 45(1), 5-32. doi: 10.1023/A:
 1010933404324
- Burrows, E. H., Bubier, J. L., Mosedale, A., Cobb, G. W., & Crill, P. M. (2005). Net
 ecosystem exchange of carbon dioxide in a temperate poor fen: a comparison
 of automated and manual chamber techniques. *Biogeochemistry*, 76(1), 21–45.
- Chen, N., Wang, A., An, J., Zhang, Y., Ji, R., Jia, Q., ... Guan, D. (2020). Modeling canopy carbon and water fluxes using a multilayered model over a temperate meadow in inner mongolia. *International Journal of Plant Production*, 14(1), 141-154. doi: 10.1007/s42106-019-00074-4
- Cover, T. M., & Thomas, J. A. (2012). *Elements of information theory*. John Wiley & Sons.
- ⁸⁷⁵ Dou, X., & Yang, Y. (2018). Comprehensive evaluation of machine learning ⁸⁷⁶ techniques for estimating the responses of carbon fluxes to climatic forces in ⁸⁷⁷ different terrestrial ecosystems. *Atmosphere*, 9(3). doi: 10.3390/atmos9030083
- ⁸⁷⁸ Dou, X., Yang, Y., & Luo, J. (2018). Estimating forest carbon fluxes using ma-⁸⁷⁹ chine learning techniques based on eddy covariance measurements. Sustainabil-⁸⁸⁰ ity, 10(1). doi: 10.3390/su10010203
- Drewry, D. T., Kumar, P., Long, S., Bernacchi, C., Liang, X. Z., & Sivapalan, M. (2010a). Ecohydrological responses of dense canopies to environmental variabil-

883	ity: 1. Interplay between vertical structure and photosynthetic pathway. Journal
884	of Geophysical Research: Biogeosciences, 115(4). doi: 10.1029/2010JG001340
885	Drewry, D. T., Kumar, P., Long, S., Bernacchi, C., Liang, X. Z., & Sivapalan, M.
886	(2010b). Ecohydrological responses of dense canopies to environmental vari-
887	ability: 2. Role of acclimation under elevated CO_2 . Journal of Geophysical
888	Research: Biogeosciences, $115(4)$, 1–22. doi: $10.1029/2010$ JG001341
889	Dunkl, I., Spring, A., Friedlingstein, P., & Brovkin, V. (2021). Process-based analysis
890	of terrestrial carbon flux predictability. Earth System Dynamics, 12(4), 1413–
891	1426. doi: 10.5194/esd-12-1413-2021
892	Dutta, D., Wang, K., Lee, E., Goodwell, A., Woo, D., Wagner, D., & Kumar, P.
893	(2017). Characterizing vegetation canopy structure using airborne remote
894	sensing data. IEEE Transactions on Geoscience and Remote Sensing, 55(2),
895	1160–1178. doi: 10.1109/TGRS.2016.2620478
896	Falge, E., Baldocchi, D., Tenhunen, J., Aubinet, M., Bakwin, P., Berbigier, P.,
897	Wofsy, S. (2002). Seasonality of ecosystem respiration and gross
898	primary production as derived from fluxnet measurements. <i>Agricultural</i>
899	and Forest Meteorology, 113(1), 53-74. (FLUXNET 2000 Synthesis) doi:
900	https://doi.org/10.1016/S0168-1923(02)00102-8
901	Farahani, M. A., Vahid, A., & Goodwell, A. (2022). Evaluating ecohydrological model
902	sensitivity to input variability with an information-theory-based approach. En -
903	tropy, 24(7). doi: $10.3390/e24070994$
904	Franzen, S. E., Farahani, M. A., & Goodwell, A. (2020). Information flows: Char-
905	acterizing precipitation-streamflow dependencies in the Colorado headwaters
906	with an information theory approach. $Water Resources Research, 56(10),$
907	e2019WR026133. doi: https://doi.org/10.1029/2019WR026133
908	Goodfellow, I., Bengio, Y., & Courville, A. (2016). <i>Deep learning</i> . MIT Press.
909	(http://www.deeplearningbook.org)
910	Goodwell, A., & Bassiouni, M. (2022). Source relationships and model structures de-
911	termine information flow paths in ecohydrologic models. Water Resources Re-
912	search, 58(9). doi: https://doi.org/10.1029/2021WR031164
913	Goodwell, A., Jiang, P., Ruddell, B. L., & Kumar, P. (2020). Debates—does
914	information theory provide a new paradigm for Earth science? Causal-
915	ity, interaction, and feedback. Water Resources Research, $56(2)$. doi:
916	https://doi.org/10.1029/2019WR024940
917	Goodwell, A., & Kumar, P. (2017). Temporal information partitioning: Characterizing
918	synergy, uniqueness, and redundancy in interacting environmental variables. Wa -
919	ter Resources Research, 5920–5942. doi: 10.1002/2016WR020218
920	Gupta, H. V., Kling, H., Yilmaz, K. K., & Martinez, G. F. (2009). Decomposi-
921	tion of the mean squared error and nse performance criteria: Implications for
922	improving hydrological modelling. Journal of Hydrology, $377(1)$, 80-91. doi:
923	https://doi.org/10.1016/j.jhydrol.2009.08.003
924	Hawkins, L. R., Bassouni, M., Anderegg, W. R. L., Venturas, M. D., Good, S. P.,
925	Kwon, H. J., Still, C. J. (2022). Comparing model representations of
926	physiological limits on transpiration at a semi-arid ponderosa pine site. Jour-
927	nal of Advances in Modeling Earth Systems, $14(11)$, e2021MS002927. doi:
928	https://doi.org/10.1029/2021MS002927
929	He, L., Li, J., Harahap, M., & Yu, Q. (2018). Scale-specific controller of carbon
930	and water exchanges over wheat field identified by ensemble empirical mode
931	decomposition. International Journal of Plant Production, $12(1)$, 43-52. doi:
932	10.1007/s42106-017-0005-8
933	Hernandez Rodriguez, L., Goodwell, A., & Kumar, P. (2023). Inside the flux footprint:
934	The role of organized land cover heterogeneity on the dynamics of observed
935	land-atmosphere exchange fluxes. Agricultural and Forest Meteorology. doi:
936	http://dx.doi.org/10.2139/ssrn.4034618
937	Hestness, J., Narang, S., Ardalani, N., Diamos, G., Jun, H., Kianinejad, H.,

938	Zhou, Y. (2017). Deep learning scaling is predictable, empirically.
939	Hochreiter, S., & Schmidhuber, J. (1997a). Long Short-Term Memory. Neural Compu-
940	tation, 9(8), 1735-1780. doi: 10.1162/neco.1997.9.8.1735
941	Hochreiter, S., & Schmidhuber, J. (1997b). Long Short-Term Memory. Neural Compu-
942	tation, 9(8), 1735-1780. doi: 10.1162/neco.1997.9.8.1735
943	Hrachowitz, M., Savenije, H., Blöschl, G., McDonnell, J., Sivapalan, M., Pomeroy,
944	J., Cudennec, C. (2013). A decade of predictions in ungauged basins
945	(pub)—a review. Hydrological Sciences Journal, 58(6), 1198-1255. doi:
946	10.1080/02626667.2013.803183
947	Huang, CW., Domec, JC., Ward, E. J., Duman, T., Manoli, G., Parolari, A. J.,
948	& Katul, G. G. (2017). The effect of plant water storage on water fluxes
949	within the coupled soil-plant system. New Phytologist, $213(3)$, 1093-1106. doi:
950	10.1111/nph.14273
951	Jensen, L., Mueller, T., Tate, K., Ross, D., Magid, J., & Nielsen, N. (1996). Soil
952	surface co2 flux as an index of soil respiration in situ: A comparison of two
953	chamber methods. Soil Biology and Biochemistry, 28(10), 1297-1306. doi:
954	https://doi.org/10.1016/S0038-0717(96)00136-8
955	Jung, M., Reichstein, M., Margolis, H. A., Cescatti, A., Richardson, A. D., Arain,
956	M. A., Williams, C. (2011). Global patterns of land-atmosphere fluxes of
957	carbon dioxide, latent heat, and sensible heat derived from eddy covariance,
958	satellite, and meteorological observations. Journal of Geophysical Research:
959	<i>Biogeosciences</i> , 116(G3). doi: https://doi.org/10.1029/2010JG001566
960	Knoben, W. J. M., Freer, J. E., & Woods, R. A. (2019). Technical note: Inher-
961	ent benchmark or not? comparing nash-sutcliffe and kling-gupta efficiency
962	scores. Hydrology and Earth System Sciences, 23(10), 4323–4331. doi:
963	10.5194/hess-23-4323-2019
964	Kratzert, F., Klotz, D., Brenner, C., Schulz, K., & Herrnegger, M. (2018). Rainfall-
965	runoff modelling using long short-term memory (lstm) networks. Hydrology and $E_{\rm rul}$ (l, $G_{\rm rul}$, $G_{\rm $
966	Earth System Sciences, $22(11)$, $6005-6022$. doi: $10.5194/ness-22-6005-2018$
967	(2010) Towarda learning universal perional and least hydrological hebring via
968	(2019). Towards learning universal, regional, and local hydrological benaviors via
969	Sciences 22(12) 5080 5110 doi: 10.5104/boss 23.5080.2010
970	$L_{0} = P_{1} V_{1} V_{1} V_{1} V_{2} V_$
971	Le, F. V. V., Kullai, F., & Diewry, D. I. (2011). Implications for the hydrologic
972	midwestern united states Proceedings of the National Academy of Sciences
973 974	108(37), 15085-15090. doi: 10.1073/pnas.1107177108
975	Leclerc, M. Y., & Foken, T. (2014). Footprints in micrometeorology and ecology
976	(Vol. 239). Springer.
977	Leroux, L., Bégué, A., Seen, D. L., Jolivot, A., & Kayitakire, F. (2017). Driving forces
978	of recent vegetation changes in the Sahel: Lessons learned from regional and
979	local level analyses. Remote Sensing of Environment, 191, 38–54.
980	Liang, J., Guo, Q., Zhang, Z., Zhang, M., Tian, P., & Zhang, L. (2020). Influence
981	of complex terrain on near-surface turbulence structures over loess plateau
982	(Vol. 11) (No. 9). doi: 10.3390/atmos11090930
983	Meng, Y., Yang, M., Liu, S., Mou, Y., Peng, C., & Zhou, X. (2021). Quantitative
984	assessment of the importance of bio-physical drivers of land cover change based
985	on a random forest method. <i>Ecological Informatics</i> , 61, 101204.
986	Meyer, H., Reudenbach, C., Wöllauer, S., & Nauss, T. (2019). Importance of spatial
987	predictor variable selection in machine learning applications-moving from data
988	reproduction to spatial prediction. <i>Ecological Modelling</i> , 411, 108815.
989	Minns, A. W., & Hall, M. J. (1996). Artificial neural networks as rainfall-
990	runon models. $Hyarological Sciences Journal, 41(3), 399-417.$ doi: 10.1080/026266666600401511
991	10.1000/02020009009491011 Moffet A M Dapolo D Boichstein M Hellinger D V Dishandson A D Down
992	monae, m. m., i apaie, D., metenstein, m., noninger, D. I., Mchardson, A. D., Daff,

993 994	A. G., Stauch, V. J. (2007). Comprehensive comparison of gap-filling techniques for eddy covariance net carbon fluxes. <i>Agricultural and Forest</i>
995	Meteorology, 147(3), 209-232.
996	Moges, E., Ruddell, B. L., Zhang, L., Driscoll, J. M., Norton, P., Perez, F., &
997	Larsen, L. G. (2022). Hydrobench: Jupyter supported reproducible hydrological
998	model benchmarking and diagnostic tool. Frontiers in Earth Science, 10. doi:
999	10.3389/feart.2022.884766
1000	Nash, J., & Sutcliffe, J. (1970). River flow forecasting through conceptual models
1001	part 1 — a discussion of principles. Journal of Hydrology, $10(3)$, 282-290. doi:
1002	https://doi.org/10.1016/0022-1694(70)90255-6
1003	Nearing, G. S., Mocko, D. M., Peters-Lidard, C. D., Kumar, S. V., & Xia, Y. (2016).
1004	Benchmarking NLDAS-2 soil moisture and evapotranspiration to separate un-
1005	certainty contributions. Journal of Hydrometeorology (2013), 160113112628008,
1006	$\begin{array}{c} \text{doi: } 10.11(5/\text{JHM}\text{-}\text{D}\text{-}15\text{-}0005.1 \\ \text{N} & \cdot & $
1007	(2020) Deer Information Theory Duraida a New Deer directory Factoria
1008	(2020). Does information Theory Provide a New Paradigm for Earth Sci-
1009	https://doi.org/10.1020/2010WD024018
1010	Nearing C C Buddell D I Clark M D Niigaan D & Datara Lidard C (2018)
1011	Reaching, G. S., Ruddell, D. L., Clark, M. P., Nijssell, D., & Peters-Lidard, C. (2016).
1012	rology 10(11) 1835 1852 doi: 10.1175/IHM D.17.0200.1
1013	Ooba M Hirana T Morami I I Hirata B l_2 Fujinuma V (2006) Com
1014	parisons of gap-filling methods for carbon flux dataset: A combination of a
1015	genetic algorithm and an artificial neural network Ecological Modelling 198(3)
1010	
1017	Papale D Black T A Carvalhais N Cescatti A Chen J Jung M
1010	Ráduly B (2015) Effect of spatial sampling from european flux tow-
1020	ers for estimating carbon and water fluxes with artificial neural networks.
1021	Journal of Geophysical Research: Biogeosciences, 120(10), 1941-1957. doi:
1022	10.1002/2015JG002997
1023	Papale, D., & Valentini, R. (2003). A new assessment of european forests carbon
1024	exchanges by eddy fluxes and artificial neural network spatialization. $Global$
1025	Change Biology, $9(4)$, 525-535. doi: 10.1046/j.1365-2486.2003.00609.x
1026	Pastorello, G., Trotta, C., Canfora, E., Chu, H., Christianson, D., Cheah, YW.,
1027	Papale, D. (2020). The FLUXNET2015 dataset and the oneflux pro-
1028	cessing pipeline for eddy covariance data. Scientific Data, $7(1)$, 225. doi:
1029	10.1038/s41597-020-0534-3
1030	Prueger, J., & Parkin, T. (2016a). Ameriflux base us-br1 brooks field site 10- ames.
1031	AmeriFlux AMP, (Dataset). doi: https://doi.org/10.17190/AMF/1246038
1032	Prueger, J., & Parkin, T. (2016b). Ameriflux base us-br3 brooks field site 11- ames.
1033	AmeriFlux AMP, (Dataset). doi: https://doi.org/10.1/190/AMF/1246039
1034	Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalnais, N.,
1035	& Prabhat. (2019). Deep learning and process understanding for data-driven Earth systems service. Nature $\mathcal{ECC}(7742)$ data to 1028 (241506 010 0012 1
1036	Earth system science. Nature, $300(7443)$. doi: 10.1038/s41380-019-0912-1
1037	Reitz, O., Grai, A., Schmidt, M., Ketzler, G., & Leuchner, M. (2021). Opscaling net
1038	and of Coonbusical Research: Biogeographics 126(2), c2020 IC005814
1039	Ruddell B I Drowry D T & Nearing C S (2010) Information Theory for
1040	Model Diagnostics: Structural Error is Indicated by Trade Off Retwoon Func
1042	tional and Predictive Performance Water Resources Research 55(8) 6534_6554
1043	doi: 10.1029/2018WR023692
1044	Ruddell, B. L., & Kumar, P. (2009a). Ecohydrologic process networks: 1. Identifica-
1045	tion. Water Resources Research, $45(3)$, 1–23. doi: 10.1029/2008WR007279
1046	Ruddell, B. L., & Kumar, P. (2009b). Ecohydrologic process networks: 1. Identifica-
1047	tion. Water Resources Research, 45(3), 1–22. doi: 10.1029/2008WR007279

1048	Runge, J., Bathiany, S., Bollt, E., Camps-Valls, G., Coumou, D., Deyle, E.,
1049	Zscheischler, J. (2019). Inferring causation from time series in Earth system sci-
1050	ences. Nature Communications, $10(1)$, 2553. doi: $10.1038/s41467-019-10105-3$
1051	Schmidhuber, J. (2015). Deep learning in neural networks: An overview. Neural Net-
1052	works, 61, 85–117. doi: 10.1016/j.neunet.2014.09.003
1053	Schreiber, T. (2000). Measuring information transfer. <i>Physical Review Letters</i> , 85(2),
1054	461. doi: 10.1103/PhysRevLett.85.461
1055	Sendrowski, A., & Passalacqua, P. (2017). Process connectivity in a naturally pro-
1056	grading river delta. Water Resources Research, $53(3)$, 1841–1863. doi: 10.1002/
1057	2016 WR019768
1058	Sendrowski, A., Sadid, K., Meselhe, E., Wagner, W., Mohrig, D., & Passalacqua, P.
1059	(2018). Transfer entropy as a tool for hydrodynamic model validation. <i>Entropy</i> ,
1060	20(1). doi: 10.3390/e20010058
1061	Shannon, C. (1948). A mathematical theory of communication. The Bell System Tech-
1062	nical Journal, $196(4)$, $519-520$. doi: $10.1016/S0016-0032(23)90506-5$
1063	Sivapalan, M. (2003). Prediction in ungauged basins: A grand challenge for theoreti-
1064	cal hydrology. Hydrological Processes, 17, 3163 - 3170. doi: 10.1002/hyp.5155
1065	Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014).
1066	Dropout: A simple way to prevent neural networks from overfitting. J. Mach.
1067	Learn. Res., $15(1)$, 1929–1958.
1068	Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learn-
1069	ing with neural networks (Vol. 27; Z. Ghahramani, M. Welling, C. Cortes,
1070	N. Lawrence, & K. Weinberger, Eds.). Curran Associates, Inc. doi:
1071	https://doi.org/10.48550/arXiv.1409.3215
1072	Suyker, A. (2022a). Ameriflux base us-nel mead - irrigated continuous maize site.
1073	AmeriFlux AMP, (Dataset). doi: https://doi.org/10.17190/AMF/1246084
1074	Suyker, A. (2022b). Ameriflux base us-ne2 mead - irrigated maize-soybean rotation
1075	site. AmeriFlux AMP, (Dataset). doi: https://doi.org/10.1/190/AMF/1240085
1076	Suyker, A. (2022c). Ameriflux base us-ne3 mead - rainfed maize-soybean rotation site.
1077	Tennent C Largen I Pollugi D Magoz E Zhang I & Ma H (2020) The
1078	utility of information flow in formulating discharge forecast models: A case
1079	study from an arid snow-dominated catchment Water Resources Research
1080	56(8) = 2019 WR024908 doi: 10.1020/2019 WR024908
1082	Tramontana C. Migliavacca M. Jung M. Beichstein M. Keenan T. F. Camps-
1002	Valls G Papale D (2020) Partitioning net carbon dioxide fluxes into
1084	photosynthesis and respiration using neural networks <i>Global Change Biology</i>
1085	26(9), 5235-5253, doi: https://doi.org/10.1111/gcb.15203
1086	Vetter, M., Churkina, G., Jung, M., Beichstein, M., Zaehle, S., Bondeau, A.,
1087	Heimann, M. (2008). Analyzing the causes and spatial pattern of the European
1088	2003 carbon flux anomaly using seven models. <i>Biogeosciences</i> , 5(2), 561–583.
1089	doi: 10.5194/bg-5-561-2008
1090	Wang, D., Ricciuto, D., Post, W., & Berry, M. W. (2011). Terrestrial ecosystem car-
1091	bon modeling. In D. Padua (Ed.), Encyclopedia of parallel computing (p. 2034-
1092	2039). Boston, MA: Springer US. doi: 10.1007/978-0-387-09766-4_395
1093	Wang, T., Brender, P., Ciais, P., Piao, S., Mahecha, M. D., Chevallier, F., Vac-
1094	cari, F. P. (2012). State-dependent errors in a land surface model across biomes
1095	inferred from eddy covariance observations on multiple timescales. <i>Ecological</i>
1096	Modelling, 246, 11-25.
1097	Welchowski, T., Maloney, K. O., Mitchell, R., & Schmid, M. (2022). Techniques to
1098	improve ecological interpretability of black-box machine learning models. Jour-
1099	nal of Agricultural, Biological and Environmental Statistics, 27(1), 175-197. doi:
1100	10.1007/s13253-021-00479-7
1101	Williams, P. L., & Beer, R. D. (2010). Nonnegative decomposition of multivariate
1102	information. arXiv preprint arXiv:1004.2515.

-31-

1102

1103	Woo, D. K., & Kumar, P. (2017). Role of micro-topographic variability on the distri-
1104	bution of inorganic soil-nitrogen age in intensively managed landscape. Water
1105	Resources Research, 53(10), 8404-8422. doi: 10.1002/2017WR021053
1106	Xiao, J., Chevallier, F., Gomez, C., Guanter, L., Hicke, J. A., Huete, A. R.,
1107	Zhang, X. (2019). Remote sensing of the terrestrial carbon cycle: A review of
1108	advances over 50 years. Remote Sensing of Environment, 233, 111383. doi:
1109	https://doi.org/10.1016/j.rse.2019.111383
1110	Xiao, J., Davis, K. J., Urban, N. M., Keller, K., & Saliendra, N. Z. (2011). Upscaling
1111	carbon fluxes from towers to the regional scale: Influence of parameter vari-
1112	ability and land cover representation on regional flux estimates. Journal of
1113	Geophysical Research: Biogeosciences, 116(G3).
1114	Xu, M., & Qi, Y. (2001). Soil-surface CO_2 efflux and its spatial and temporal
1115	variations in a young ponderosa pine plantation in Northern California. Global
1116	Change Biology, 7(6), 667-677. doi: https://doi.org/10.1046/j.1354-1013.2001
1117	.00435.x
1118	Yan, Q., Le, P. V. V., Woo, D. K., Hou, T., Filley, T., & Kumar, P. (2019). Three-
1119	dimensional modeling of the coevolution of landscape and soil organic carbon.
1120	Water Resources Research, 55(2), 1218-1241. doi: 10.1029/2018WR023634
1121	Zhou, Q., Fellows, A., Flerchinger, G. N., & Flores, A. N. (2019). Examining inter-
1122	actions between and among predictors of net ecosystem exchange: A machine
1123	learning approach in a semi-arid landscape. Scientific Reports, $9(1)$, 2222. doi:
1124	10.1038/s41598-019-38639-y
1125	Zhou, X., Wang, X., Tong, L., Zhang, H., Lu, F., Zheng, F., Ouyang, Z. (2012).
1126	Soil warming effect on net ecosystem exchange of carbon dioxide during the
1127	transition from winter carbon source to spring carbon sink in a temperate
1128	urban lawn. Journal of Environmental Sciences, 24(12), 2104-2112. doi:
1129	https://doi.org/10.1016/S1001-0742(11)61057-7

-32-

Supporting Information for "Causal Drivers of Land-Atmosphere Carbon Fluxes from Machine Learning Models and Data"

Mozhgan A. Farahani¹, Allison E. Goodwell^{1,2}

 $^1 \mathrm{University}$ of Colorado Denver, Department of Civil Engineering

 $^2\mathrm{Prairie}$ Research Institute, University of Illinois at Urbana-Champaign

Contents of this file

- 1. Text S1 to S4
- 2. Figures S1 to S18
- 3. Tables S1

S1. Data pre-processing

The data pre-processing stage was a crucial step in our study, ensuring the reliability and accuracy of our analysis. This process involved several steps:

1.1. Quality Control

Firstly, we applied quality control measures to all datasets. This involved checking for any inconsistencies, errors, or outliers in the data that could potentially skew our results. We used a combination of automated checks and manual review to ensure the integrity

Corresponding author: A. E. Goodwell, Department of Civil Engineering, University of Colorado Denver, USA., (Alison.goodwell@ucdenver.edu) of our data. Automated checks included algorithms to detect statistical anomalies, while manual review involved visual inspection of the data and cross-checking with source documentation.

1.2. Handling Missing Values

In some datasets, we encountered missing values. To handle these, we used time series imputation methods. The choice of imputation method was dependent on the distribution of the data. For normally distributed data, we used mean imputation. This technique replaces the missing values with the average of the available data for that variable, thus capitalizing on the characteristic symmetric nature of the distribution. Specifically, the variables Fc, SWC, Ta, TS, and Pa were treated using mean imputation. Conversely, for those variables presenting skewed distributions or characterized by extreme outliers, median imputation was employed. The median, being the middle value of a dataset, is less sensitive to outliers and provides a more robust measure of central tendency for skewed distributions. The variables WS, P, NETRAD, PPFD and RH were imputed using this method. Through these imputation strategies, we ensured that the integrity of the data distribution was upheld, while concurrently addressing the gaps in our dataset.

Moreover, to address significant missing values in the *PPFD* variable at Br1 and Br3 sites, we employed a linear regression imputation technique using *NETRAD* values as predictors. We first used those part of datasets where *PPFD* and *NETRAD* were concurrently present, using them as training data for individual linear regression models. Once trained, these models were used to predict missing *PPFD* values based on available *NETRAD* values, thus leveraging their linear relationship for accurate imputation.

1.3. Normalization

To ensure efficient learning and to prevent any one variable from dominating others due to scale differences, we normalized all input variables and the output (Fc) data. Specifically, we utilized the "MinMax" scaling technique, where the minimum of feature is made equal to zero and the maximum of feature equal to one. In this method, every feature value is transformed to fall within the range [0,1]. It scales the values to the specific value range without changing the shape of the original distribution. This approach entails subtracting the minimum value of the feature and then dividing by the range of that feature, resulting in a dataset where the minimum and maximum feature values are normalized to lie between 0 and 1. This procedure not only enhances the efficiency of learning algorithms but also aids in preventing potential numerical stability issues.

1.4. Retransformation

The output of all machine learning models was retransformed using the normalization parameters to obtain the final Fc prediction in the original scale. This step is crucial for interpreting the results in the context of the original data.

It's important to note that while these pre-processing steps greatly enhance the quality and usability of the data, they are based on certain assumptions and can introduce some level of uncertainty. However, we applied these methods systematically and transparently to minimize potential biases and ensure the reliability of our results. The full suite of variables used in this study, along with their descriptions and units, is outlined in Table ?? in the main manuscript.

S2. Information Decomposition We use information decomposition to analyze causal interactions in which two sources provide information to a target variable, which could be

an observation or a model output. In a system where two sources share information from X and Y with a target Z, the total information quantity, I(X,Y;Z), can be partitioned into synergistic (S), unique (U), and redundant (R) components. Any existing IT-based measure can also be defined in terms of combinations of R, U, and S (Figure S1). For example, this partitioning of information implies that the mutual information between the target and each source is the sum of the redundancy and the unique information from the source, i.e. $I(X;Z) = U_{X|Y} + R_{X,Y}$ (Figure S1a). Meanwhile, conditional mutual information, which includes transfer entropy as a special case, is the sum of unique and synergistic components, i.e. $I(X;Z|Y) = U_{X|Y} + S_{X,Y}$ (Figure S1b). Finally, the interaction information, which is symmetric between all three variables, is equivalent to $S_{X,Y} - R_{X,Y}$ (Goodwell & Kumar, 2017, 2015), such that positive or negative interaction information indicates whether synergy or redundancy is dominant (Figure S1c). To simplify notation hereafter, we omit subscripts such that $S_{X,Y} = S$ and $R_{X,Y} = R$ given a particular definition of sources and target. We similarly simplify unique information components to $U_{X|Y} = U_X$ and $U_{Y|X} = U_Y$.

While information decomposition is a useful concept, information theory does not provide formulas to directly determine these quantities. Several studies (Barrett, 2015; Williams & Beer, 2010) defined redundancy measures as the mutual information that the weakest source provides to the target, forcing one unique component to equal zero. Goodwell and Kumar considered that this is actually a maximum bound for redundancy, and applied a "rescaled" redundancy measure in which redundancy is scaled between the minimum and maximum bounds that are defined by information theory. The maximum bound is the minimum mutual information that either source provides to the target,

 $R_{max} = min[I(X;Z), I(Y;Z)]$. The minimum bound is zero for cases where the interaction information is positive or $S_R > 0$, i.e. I(X,Y;Z) > I(X;Z) + I(Y;Z). otherwise, if S - R < 0, the minimum bound for redundancy is the negative interaction information, in order for synergy to be non-negative. This leads to a definition of the minimum R as $R_{min} = max[0, I(X;Z) + I(Y;Z) - I(X,Y;Z)]$. We then scale redundancy between these bounds based on the normalized information between the source variables:

$$I_s = \frac{I(X;Y)}{\min[H(X),H(Y)]} \tag{1}$$

$$R_s = R_{\min} + I_s (R_{\max} - R_{\min})$$

In general, this definition causes highly correlated sources to be maximally redundant with each other, while independent sources are minimally redundant. A definition for redundancy enables the computation of the other information decomposition components, S, U_X , and U_Y .

S3. Statistical Significance

We compute statistical significance of observed or modeled information theoretic measures using a shuffled surrogates approach. We define a critical value of total information as follows:

$$I_{crit} = I_{suff, mean} + 3 \times I_{suff, stdev}$$
(2)

where $I_{suff,mean}$ and $I_{suff,stdev}$ are the mean and standard deviation of 100 information measures computed with randomly shuffled source data. For example, if the $I(Ta, Ts; Fc) < I_{crit}$, we set all information components to zero and do not do fur-

ther information partitioning. Meanwhile, if I(Ta, Ts; Fc) is statistically significant but I(Ta; Fc|Ts) is not (according to the same shuffled surrogate method), we set the unique component from Ta and the synergistic component to zero, since $I(Ta; Fc|Ts) = U_{Ta} + S$. Then, we define R as I(Ta; Fc), since $I(Ta; Fc) = U_{Ta} + R$, and U_{Ts} is computed as $U_{Ts} = I(Ta; Ts; Fc) - R$. For a case where I(Ts; Fc|Ta) is not statistically significant, we apply a similar process. Finally, if neither conditional term of I(Ta; Fc|Ts) or I(Ts; Fc|Ta) is statistically significant would indicate that the only information component is redundancy. However, we defined that this case never occurs based on our study year period.

S4. Functional Performance

We calculated the individual level (Figures S3, S4, S5, S6, S7) and pairwise level (Figures S8, S9, S10, S11, S12, S14, S15, S16, S17, S18) of functional performance at Ne2, Ne3, Br1, Br3 and GC sites. These sites show similar patterns in mutual information as site Ne1 which presented in the main manuscript. We also find similar patterns in pairwise functional performance, specifically the overestimation of U at the expense of S and overestimation of R for correlated source pairs. However, we find that regionally trained models (Figures S13-S18) diminish some of the issues observed in the localized models (Figures S8-S12). The regional model also corrects the balance between synergy and unique contributions, leading to a more accurate representation of how these variables interact. This trend is especially evident in the LSTM model, which demonstrates enhanced functional performance under regional training.

References

Barrett, A. B. (2015). Exploration of synergistic and redundant information sharing in

static and dynamical Gaussian systems. Physical Review E, 91(5). doi: 10.1103/ PhysRevE.91.052802

- Goodwell, A., & Kumar, P. (2015). Information theoretic measures to infer feedback dynamics in coupled logistic networks. *Entropy*, 17(11), 7468–7492. doi: 10.3390/ e17117468
- Goodwell, A., & Kumar, P. (2017). Temporal information partitioning: Characterizing synergy, uniqueness, and redundancy in interacting environmental variables. Water Resources Research, 5920–5942. doi: 10.1002/2016WR020218
- Williams, P. L., & Beer, R. D. (2010). Nonnegative decomposition of multivariate information. arXiv preprint arXiv:1004.2515.



Figure S1. Illustration of information theory metrics. (a) Mutual information I(X; Z) is the reduction in uncertainty about Z given knowledge of X. (b) Conditional mutual information I(X; Z|Y) is the reduction in uncertainty about Z given knowledge of X, beyond information already provided by Y. (c) Multi-variate mutual information I(X, Y; Z) is the total reduction in uncertainty about Z given knowledge of X and Y together, and is composed of four non-negative components of R, U_X , U_Y , and S.

Attribute	Description/Value
Model Type	Multiple Linear Regression (MLR)
Method	Ordinary Least Squares (OLS)
Implementation	"statsmodels" package in Python
Model Type	Random Forest (RF)
Trees in the Forest	100 (n-estimators)
Max Features	Square root of total features
Structure	Ensemble of Decision Trees
Implementation	"scikit-learn" package in Python
Model Type	Long Short Term Memory Model (LSTM)
Number of LSTM Layers	2
Number of Hidden Units per Layer	9
Dropout Layers	Between LSTM layers
Final Layer Type	Regression $(1 \text{ unit for } Fc)$
Sequence Length	12 time steps (half a diurnal cycle)
Batch Size	128
Loss Function	Mean Squared Error (MSE)
Implementation	"torch" package in Python

 Table S1.
 Summary of Machine Learning Model Architecture

X - 8



Figure S2. Averaged monthly values of driving variables (air temperature (Ta), relative humidity (RH), precipitation (P), soil temperature (TS), photosynthetic photon flux density (PPFD), net radiation (NETRAD), wind speed (WS), atmospheric pressure (Pa), soil water content (SWC)) and target variable (Fc) over the study years corresponded to different sites (Ne1, Ne2, Ne3, Br1, Br3, GC). Each site is represented by a unique color.



Figure S3. (a) Normalized mutual information (I_n) and (b) Individual source level of functional performance $(A_{f,MI})$ of three different models - Multiple Linear Regression (MLR), Random Forest (RF), and Long Short-Term Memory (LSTM) - under two training experiences, local and regional, at Ne2 site. Each variable is ranked based on the average observed MI across all sites. Observation values are represented with a black dot linked by a dashed line.



Figure S4. (a) Normalized mutual information (I_n) and (b) Individual source level of functional performance $(A_{f,MI})$ of three different models - Multiple Linear Regression (MLR), Random Forest (RF), and Long Short-Term Memory (LSTM) - under two training experiences, local and regional, at Ne3 site. Each variable is ranked based on the average observed MI across all sites. Observation values are represented with a black dot linked by a dashed line.



Figure S5. (a) Normalized mutual information (I_n) and (b) Individual source level of functional performance $(A_{f,MI})$ of three different models - Multiple Linear Regression (MLR), Random Forest (RF), and Long Short-Term Memory (LSTM) - under two training experiences, local and regional, at Br1 site. Each variable is ranked based on the average observed MI across all sites. Observation values are represented with a black dot linked by a dashed line.



Figure S6. (a) Normalized mutual information (I_n) and (b) Individual source level of functional performance $(A_{f,MI})$ of three different models - Multiple Linear Regression (MLR), Random Forest (RF), and Long Short-Term Memory (LSTM) - under two training experiences, local and regional, at Br3 site. Each variable is ranked based on the average observed MI across all sites. Observation values are represented with a black dot linked by a dashed line.



Figure S7. (a) Normalized mutual information (I_n) and (b) Individual source level of functional performance $(A_{f,MI})$ of three different models - Multiple Linear Regression (MLR), Random Forest (RF), and Long Short-Term Memory (LSTM) - under two training experiences, local and regional, at GC site. Each variable is ranked based on the average observed MI across all sites. Observation values are represented with a black dot linked by a dashed line.



Figure S8. Observed pairwise (a) synergistic $(S_{i,j})$, (b) redundancy $(R_{i,j})$, and (c) uniqueness $(U_{i,j})$ information flow at Ne2 site. Pairwise functional performance of three models under local training experience at Ne2 site. The heat-map represents the relative difference in information decomposition partitioning measures $(A_{f,S_{i,j}}, A_{f,R_{i,j}}, \text{ and } A_{f,U_{i,j}})$ between modeled and observed data for each pair of forcing variables. Positive values (green) in (d)-(l) indicate that the model overestimates the information type, while negative values (red) indicate underestimations. September 6, 2023, 8:01pm



Figure S9. Observed pairwise (a) synergistic $(S_{i,j})$, (b) redundancy $(R_{i,j})$, and (c) uniqueness $(U_{i,j})$ information flow at Ne3 site. Pairwise functional performance of three models under local training experience at Ne3 site. The heat-map represents the relative difference in information decomposition partitioning measures $(A_{f,S_{i,j}}, A_{f,R_{i,j}}, \text{ and } A_{f,U_{i,j}})$ between modeled and observed data for each pair of forcing variables. Positive values (green) in (d)-(l) indicate that the model overestimates the information type, while negative values (red) indicate underestimations. September 6, 2023, 8:01pm



Figure S10. Observed pairwise (a) synergistic $(S_{i,j})$, (b) redundancy $(R_{i,j})$, and (c) uniqueness $(U_{i,j})$ information flow at Br1 site. Pairwise functional performance of three models under local training experience at Br1 site. The heat-map represents the relative difference in information decomposition partitioning measures $(A_{f,S_{i,j}}, A_{f,R_{i,j}}, \text{ and } A_{f,U_{i,j}})$ between modeled and observed data for each pair of forcing variables. Positive values (green) in (d)-(l) indicate that the model overestimates the information type, while negative values (red) indicate underestimations. September 6, 2023, 8:01pm



Figure S11. Observed pairwise (a) synergistic $(S_{i,j})$, (b) redundancy $(R_{i,j})$, and (c) uniqueness $(U_{i,j})$ information flow at Br3 site. Pairwise functional performance of three models under local training experience at Br3 site. The heat-map represents the relative difference in information decomposition partitioning measures $(A_{f,S_{i,j}}, A_{f,R_{i,j}}, \text{ and } A_{f,U_{i,j}})$ between modeled and observed data for each pair of forcing variables. Positive values (green) in (d)-(l) indicate that the model overestimates the information type, while negative values (red) indicate underestimations. September 6, 2023, 8:01pm



Figure S12. Observed pairwise (a) synergistic $(S_{i,j})$, (b) redundancy $(R_{i,j})$, and (c) uniqueness $(U_{i,j})$ information flow at GC site. Pairwise functional performance of three models under local training experience at GC site. The heat-map represents the relative difference in information decomposition partitioning measures $(A_{f,S_{i,j}}, A_{f,R_{i,j}}, \text{ and } A_{f,U_{i,j}})$ between modeled and observed data for each pair of forcing variables. Positive values (green) in (d)-(l) indicate that the model overestimates the information type, while negative values (red) indicate underestimations. September 6, 2023, 8:01pm



Figure S13. Pairwise functional performance of three models under regional training experience at Ne1 site. The heat-map represents the relative difference in information decomposition partitioning measures $(A_{f,S_{i,j}}, A_{f,R_{i,j}}, \text{ and } A_{f,U_{i,j}})$ between modeled and observed data for each pair of forcing variables. Positive values (green) indicate that the model overestimates the information type, while negative values (red) indicate underestimations.



Figure S14. Pairwise functional performance of three models under regional training experience at Ne2 site. The heat-map represents the relative difference in information decomposition partitioning measures $(A_{f,S_{i,j}}, A_{f,R_{i,j}}, \text{ and } A_{f,U_{i,j}})$ between modeled and observed data for each pair of forcing variables. Positive values (green) indicate that the model overestimates the information type, while negative values (red) indicate underestimations.



Figure S15. Pairwise functional performance of three models under regional training experience at Ne3 site. The heat-map represents the relative difference in information decomposition partitioning measures $(A_{f,S_{i,j}}, A_{f,R_{i,j}}, \text{ and } A_{f,U_{i,j}})$ between modeled and observed data for each pair of forcing variables. Positive values (green) indicate that the model overestimates the information type, while negative values (red) indicate underestimations.



Figure S16. Pairwise functional performance of three models under regional training experience at Br1 site. The heat-map represents the relative difference in information decomposition partitioning measures $(A_{f,S_{i,j}}, A_{f,R_{i,j}}, \text{ and } A_{f,U_{i,j}})$ between modeled and observed data for each pair of forcing variables. Positive values (green) indicate that the model overestimates the information type, while negative values (red) indicate underestimations.



Figure S17. Pairwise functional performance of three models under regional training experience at Br3 site. The heat-map represents the relative difference in information decomposition partitioning measures $(A_{f,S_{i,j}}, A_{f,R_{i,j}}, \text{ and } A_{f,U_{i,j}})$ between modeled and observed data for each pair of forcing variables. Positive values (green) indicate that the model overestimates the information type, while negative values (red) indicate underestimations.



Figure S18. Pairwise functional performance of three models under regional training experience at GC site. The heat-map represents the relative difference in information decomposition partitioning measures $(A_{f,S_{i,j}}, A_{f,R_{i,j}}, \text{ and } A_{f,U_{i,j}})$ between modeled and observed data for each pair of forcing variables. Positive values (green) indicate that the model overestimates the information type, while negative values (red) indicate underestimations.