

Ward Adam Scott (Orcid ID: 0000-0002-6376-0061)  
 Drummond Jennifer (Orcid ID: 0000-0002-6501-7618)  
 Hannah David M. (Orcid ID: 0000-0003-1714-1240)  
 Krause Stefan (Orcid ID: 0000-0003-2521-2248)  
 Zarnetske Jay (Orcid ID: 0000-0001-7194-5245)

## Advancing river corridor science beyond disciplinary boundaries with an inductive approach to catalyze hypothesis generation

*Submitted for publication in Hydrological Processes – Special issue – Data Science Applications in Hydrology*

### Authors:

Adam S. Ward<sup>1</sup>, Aaron Packman<sup>2</sup>, Susana Bernal<sup>3</sup>, Nicolai Brekenfeld<sup>4</sup>, Jen Drummond<sup>4</sup>, Emily Graham<sup>5</sup>, David M. Hannah<sup>4</sup>, Megan Klaar<sup>6</sup>, Stefan Krause<sup>4</sup>, Marie Kurz<sup>7</sup>, Angang Li<sup>2</sup>, Anna Lupon<sup>3</sup>, Feng Mao<sup>8</sup>, M. Eugènia Martí Roca<sup>3</sup>, Valerie Ouellet<sup>4</sup>, Todd V. Royer<sup>1</sup>, James C. Stegen<sup>5</sup>, Jay P. Zarnetske<sup>9</sup>

<sup>1</sup> O'Neill School of Public and Environmental Affairs, Indiana University, Bloomington, Indiana, USA

<sup>2</sup> Department of Civil and Environmental Engineering, Northwestern University, Evanston, Illinois, USA

<sup>3</sup> Integrative Freshwater Ecology Group, Centre for Advanced Studies of Blanes (CEAB-CSIC), Blanes, Spain

<sup>4</sup> School of Geography, Earth & Environmental Sciences, University of Birmingham, Edgbaston, Birmingham, B15 2TT, UK

<sup>5</sup> Earth and Biological Sciences Directorate, Pacific Northwest National Laboratory, Richland, Washington, USA

<sup>6</sup> School of Geography, School of Earth and Environment, University of Leeds, Woodhouse, Leeds LS2 9JT, United Kingdom

<sup>7</sup> The Academy of Natural Sciences of Drexel University, Philadelphia, Pennsylvania, USA

<sup>8</sup> School of Earth and Environmental Sciences, Cardiff University, Building, Park Place, Cardiff, CF10 3AT, United Kingdom

<sup>9</sup> Department of Earth and Environmental Sciences, Michigan State University, East Lansing, Michigan, USA

### Corresponding author:

Adam S. Ward  
 O'Neill School of Public and Environmental Affairs  
 Indiana University  
 418 MSB-II  
 Bloomington, IN 47405

Email: adamward@indiana.edu  
 Phone: 812-865-4820

**Running head:** Inductive hypothesis generation using data science

**Key words:** river corridor, stream corridor, machine learning, inductive, scientific method

This article has been accepted for publication and undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process which may lead to differences between this version and the [Version of Record](#). Please cite this article as doi: [10.1002/hyp.14540](https://doi.org/10.1002/hyp.14540)

## Abstract

A unified conceptual framework for river corridors requires synthesis of diverse site-, method- and discipline-specific findings. The river research community has developed a substantial body of observations and process-specific interpretations, but we are still lacking a comprehensive model to distill this knowledge into fundamental transferable concepts. We confront the challenge of how a discipline classically organized around the deductive model of systematically collecting of site-, scale-, and mechanism-specific observations begins the process of synthesis. Machine learning is particularly well-suited to inductive generation of hypotheses. In this study, we prototype an inductive approach to holistic synthesis of river corridor observations, using support vector machine regression to identify potential couplings or feedbacks that would not necessarily arise from classical approaches. This approach generated 672 relationships linking a suite of 157 variables each measured at 62 locations in a 5<sup>th</sup> order river network. Eighty four percent of these relationships have not been previously investigated, and representing potential (hypothetical) process connections. We document relationships consistent with current understanding including hydrologic exchange processes, microbial ecology, and the River Continuum Concept, supporting that the approach can identify meaningful relationships in the data. Moreover, we highlight examples of two novel research questions that stem from interpretation of inductively-generated relationships. This study demonstrates the implementation of machine learning to sieve complex data sets and identify a small set of candidate relationships that warrant further study, including data types not commonly measured together. This structured approach complements traditional modes of inquiry, which are often limited by disciplinary perspectives and favor the careful pursuit of parsimony. Finally, we emphasize that this approach should be viewed as a complement to, rather than in place of, more traditional, deductive approaches to scientific discovery.

## 1. Introduction

A paradigm change is required to advance our conceptualization of the river corridor beyond site-, scale-, and mechanism-specific findings towards understanding river corridors as complex, dynamic systems responding to external forcing (Turnbull et al., 2018). While decades of study have yielded descriptions of many individual process controls, we have yet to assemble this ensemble of process dynamics across space and time to create a comprehensive understanding of the structure and function of river corridors. Here and throughout we use the term ‘dynamics’ to refer to the network of couplings and feedbacks internal to a study system that stimulate mechanisms, yielding observable fluxes or state variables (consistent Stegen et al., 2018), as opposed to more narrowly describing temporal variability. Most river corridor studies focus on a specific location, scale, or disciplinary perspective, and consequently investigate a limited set of measurements (Turnbull et al., 2018; Ward, 2015; Ward & Packman, 2019). Consequently, we have accumulated a substantial body of observations and process-specific interpretations, but we are lacking a comprehensive model to distill this knowledge into general and transferable concepts. At present, few - if any - conceptual models account for the hierarchical, multi-scale, coupled physical-chemical-biological process dynamics that give rise to the observed spatio-temporal patterns of river corridor services and functions. A new approach is needed for conceptualizing the multi-scale and multi-rate interactions that span disciplines and govern river corridors, from deep time geological processes shaping landscape uplift and evolution to contemporary rapid dynamics of microbial gene expression to future responses in suspended solid transport following fire, and every physical-chemical-biological process in between.

River corridors have classically been studied by a host of disciplines, each with primary interest in individual processes or functions (Ward, 2015). Consequently, techniques for river research are not standardized across disciplines, relevant metadata have not been specified, and common variables needed to synthesize findings across sites are not defined (Ward, 2015; Ward & Packman, 2019). Thus, the core challenges facing river corridor scientists today are (a) developing theory to overcome our limited ability to observe the full spatio-temporal complexity of river corridors (Li et al., 2021), (b) organizing river corridor science in a way that is explicitly integrative as opposed to disciplinary, and (c) facilitating communication and idea generation across disciplines. One way to address these needs is to expand beyond the traditional, deductive

Accepted Article

approach to science, which bases measurements on a highly targeted set of causal mechanisms to be tested at a limited range of locations and scales. With the emergence of new experimental and data science techniques, the time has come to expand existing conceptual models for river corridors via approaches that generate more integrative knowledge commensurate with the reality of river corridors as complex dynamic systems. We posit that unified understanding must be derived from a combination of *deductive* science and *inductive* approaches that identify process interactions and couplings that emerge from the data themselves. We suggest that river corridor science can benefit from inductive approaches that generate hypotheses and eventually theories from empirical studies, an approach successfully applied in other disciplines (Martin & Turner, 1986; Strauss & Corbin, 1994; e.g., Turnbull et al., 2018).

A unifying framework is required to organize and synthesize our understanding of river corridors and advance scientific understanding of the drivers and controls of their functioning. Stegen et al. (2018) propose one such model for microbial ecology, where the resultant ecosystem functions and services are explained by the relationships linking internal dynamics, external forcing, and historical contingencies. The principles of Stegen et al.'s conceptual framework are similar to other existing conceptualizations of river corridors that have been developed by other disciplines. First, external forcing describes the role of factors extrinsic to the river corridor that shape its structure and function. For river corridors, this primarily means the larger spatial scale and longer temporal scale elements that are functionally decoupled (e.g., static or slowly-varying) relative to a process of interest. Studies with data collection spanning gradients in land use, geologic setting, climate, network position, or other factors that are considered to be extrinsic typically use geospatial and statistical approaches to describe patterns and trends (e.g., McGuire et al., 2014), while variation around spatially structured trends is often interpreted as random noise from structural heterogeneity and/or unstudied, smaller-scale processes (Abbott et al., 2018). Next, internal dynamics are the interacting processes within the river corridor that give rise to observed functions of interest at a given location. Conceptual models based on this approach to river corridor science include hot spots and hot moments (Krause et al., 2011, 2017; Wallis et al., 2020), control points (Bernhardt et al., 2017), and patch dynamics (Pringle et al., 1988). River corridor dynamics are commonly studied through detailed observations at a relatively limited spatial scale, which is restricted in an attempt to characterize local feedbacks

between mechanisms. These approaches often lack sufficient spatial resolution to enable confident application of geostatistical approaches, and may not reliably support assessments of system dynamics (e.g., Lee-Cullin et al., 2018). Longer-term dynamics are often considered as historical contingencies: the biotic and abiotic histories or antecedent conditions that lead to the present characteristics of the river corridor and affect its response to future perturbations. Examples of river corridor studies that incorporate historical contingencies include perturbation-response dynamics, commonly associated with floods (Czuba et al., 2019; Wu et al., 2018), droughts (Boulton et al., 2004; Wood et al., 2010), or restoration activities (Rana et al., 2017; Smidt et al., 2015), and large-scale historical perturbations such as land development (Liébault & Piégay, 2002; Walling & Fang, 2003; Wohl, 2005), river regulation (Gregory, 2006), and contamination (Byrne et al., 2012; Santschi et al., 2001). Such studies often involve little to no replication and may be biased towards response variables that change rapidly relative to processes that are quasi-steady over the timeframe of a given experiment.

While external forcing, internal dynamics, and historical contingencies have each been studied in their own right, recent studies are beginning to integrate these concepts into holistic understanding of river corridors. For example, Wisnoski and Lennon (2021) explicitly linked localized heterogeneity to systematic spatial patterns along the network, revealing that the local microbial assemblage in headwaters streams was controlled by local physical and chemical conditions, but these local controls gave way to systemic organization from headwaters to larger downstream rivers as the spatial scale of study increased. Such explicit consideration of local and network scales is rare and still does not address historical contingencies. However, if done more often and expanded to consider historical contingencies as a context for each replicate, this type of systematic approach would allow assessment of the transition in dominant controls from local heterogeneity (a reflection of internal dynamics) to larger-scale spatial organization (a reflection of external drivers), the specific mechanisms of this transition, and the scale at which the transition occurs, and how future multi-scale dynamics may depend on antecedent conditions (a reflection of historical contingencies). Studies that have explicitly considered local spatiotemporal dynamics as part of long-term system-wide functions have found strong relationships between large-scale system structure, internal dynamics, and long-term emergent outcomes in flow, sediment transport, and biogeochemistry (e.g., Fisher et al., 1998; Harvey &

Gooseff, 2015; Krause et al., 2017; Pinay et al., 2015). The success of these studies demonstrates our ability to identify a core set of transferable and scalable processes that govern river system dynamics and unify seemingly disparate observations into holistic understanding of river corridor services and functions.

Here we use objective data-oriented approaches to confront the challenge of how a discipline organized around the classic deductive model of site-, scale-, and mechanism-specific observations can systematically link the resulting fragmented information into system-level understanding. Our aim is to identify couplings that span scales and disciplinary expertise in absence of pre-existing conceptual models that would traditionally serve as the source of hypotheses for deductive testing. We propose an inductive approach to data synthesis, serving as a basis for the unconstrained generation of new and potentially unexpected relationships, each of which may be explained by hypotheses that could subsequently be tested. To this end, we analyze a novel large data set for a 5<sup>th</sup> order river basin (Ward, Zarnetske, et al., 2019) using inductive approaches to generate a network of relationships that span traditional disciplinary boundaries. The data set contains 157 variables with nearly 25,000 possible pairwise relationships, making it infeasible to explore each potential relationship through the lens of deductive inquiry. Further, the large degree of covariation in environmental conditions may obscure underlying causal mechanisms, making it difficult to determine unique process relationships and their controls. Thus, we pilot a machine learning approach that sieves and categorizes information to identify non-obvious relationships that merit subsequent investigation. We envision the apparent relationships generated by our approach as a suite of observations around which hypotheses could be generated and subsequently tested with more traditional approaches. In this way, we complement traditional approaches by highlighting observations that may warrant hypotheses to be spun that explain causal pathways that novel, interdisciplinary, and trans-scale to explain the apparent relationships. This allows us to synthesize complex, multi-scale observations independent of any pre-conceived conceptual models and uncover novel and exciting information about the structure and function of river corridors. We critically evaluate the resultant relationships relative to existing knowledge, and provide two examples of how these novel insights may motivate future research questions that inform a synthesis approach to understanding of river corridors.

## **2. Methods**

### ***2.1 Data description and organization***

#### ***2.1.1 Field site and synoptic campaign***

The H.J. Andrews Experimental forest (Western Cascades, Oregon, USA) is a 6,400 ha basin that is primarily covered in old-growth and second growth forest and drained by a 5<sup>th</sup> order river. The physical characteristics of the basin are well-described elsewhere (Deligne et al., 2017; Dyrness, 1969; Jefferson et al., 2004; Swanson & James, 1975; Swanson & Jones, 2002). A synoptic sampling campaign including detailed characterization of physical, chemical, and biological characteristics and processes in the river corridor at 62 sites across stream orders 1-5 was conducted by Ward et al. (2019), which forms the basis of our study data set. These data are the most uniform, comprehensive, and multi-scale available – to our knowledge – and, as such, are uniquely useful for assessment of relationships spanning scales and disciplines. Notably these data represent a spatial synoptic sampling design (i.e., a snapshot in time), meaning their analysis will necessarily highlight apparent spatial patterns but cannot capture the temporal dynamics of the system. Indeed, river corridors will have processes operation spanning orders of magnitude in temporal scale (Ward and Packman, 2019). Consequently, our approach will not capture temporal couplings between relationships, and we are combining relatively dynamic variables (e.g., water temperature) and relatively static variables (e.g., surficial geology) into a single analysis. Approaches with comparable coverage occurring through seasonal, storm, and/or diurnal fluctuations would enable a related assessment of temporal dynamics and the persistence of relationships through natural variation.

#### ***2.1.2 Data reduction***

Starting from this data set, we reduced the full suite of variables from Ward et al. (2019) to a subset we considered to be most representative summary of the data set. For example, we omitted identification of individual species and life-stages from macroinvertebrate data in favor of summary indices, and similarly reduced the 10,000+ individual organic molecules identified in the data set (i.e., metabolomics, the profiling of individual organic compounds within each sample) to a suite of summary indices. In this process, we discussed traditional disciplinary approaches to the study of river corridors, and ultimately organized the variables into 7

subgroups representing distinct study domains that jointly characterize the structure, function, and dynamics of the river corridor and consistent with the design of the field campaign. These subgroups were: geologic setting (GEO), physical chemistry (PCHEM), bulk DOM characterization (DOM), dissolved nutrients (NUTS), solute tracers (TRACER), metabolomics (ICR), and macroinvertebrates (MACRO). A complete list of variables, subgroups, and summary findings for each variable is presented in Table S1). The reduced data set totaled 157 unique variables across the seven disciplinary subgroups and is the basis for all subsequent analysis in this study.

## ***2.2 Principal components analysis***

To identify major axes of (co)variation among measured variables, we performed a series of principal component analyses (PCAs) using the rotated PCA approach. Independent PCAs were performed first on the entire data set (all 157 variables) and subsequently on variables within each subgroup. For each PCA, we focused on results from the first two components (PC1 and PC2). We identified the most influential variables from each principal component as those with loadings greater than 0.6 or less than -0.6 (hereafter ‘influential variables’) and interpreted the variables aligned with each PC to describe the major axes of variation when possible.

## ***2.3 Spatial structure of individual variables***

For each variable, we tested for spatial structure throughout the network by assessing the change in variance as a function of distance between flow connected points, (i.e., a semivariogram; Ver Hoef et al., 2006; Isaak et al., 2014; McGuire et al., 2014). This analysis identifies variables for which variance is spatially uniform (i.e., no change in variance as a function of distance), increases linearly (i.e., variance grows with distance), or that plateaus at some distance (a scale cutoff). A uniform relationship indicates no structure (hereafter, unstructured variable), while both linear and plateau relationships demonstrate spatial structure (hereafter, structured variable). The linear models were only considered significant if the estimate of the slope was significantly different from zero based on the 95% confidence interval for a linear model fit. The squared differences were normalized (squared difference subtracted from the mean, followed by division of the difference by the standard deviation) and binned (bin size of 30) before being fitted. An exponential semivariogram function was considered for cases that exhibited scale cutoffs:



$$|y = a + be^{\left(\frac{-x}{c}\right)},$$

with the `nls()` function in R Studio. The nugget, sill and range are given by  $a$ ,  $a+b$  and  $3 \times c$ , respectively. Exponential semivariogram models were only considered significant if the estimates of the parameters  $b$  and  $c$  were significantly different from zero, based on zero not being within the 95% confidence interval for the parameters.

#### ***2.4 Support vector machine regression***

To derive a network of relationships among pairs of variables in the data set, and ultimately identify the interactions within the network, we constructed two sets of support vector machine regression (SVMR) models. Each model predicted an individual dependent variable using a suite of independent variables. The model used forward feature selection with leave-one-out cross-validation. Forward selection stopped adding additional independent variables when the coefficient of determination failed to improve when an additional variable was included to limit overfitting by the model. The evaluation of each potential independent variable to add to the model was based on leave-one-out cross validation, where all possible permutations of training on all but one data point to predict the withheld data point were considered. The SVMR improvement summed across the ensemble of 62 models per variable was considered as the basis to add a variable to the feature set, and the process proceeded iteratively until adding independent variables failed to improve model fit. Gaussian kernels were used for all variables, and variables were normalized for analysis. For each SVMR we recorded the order in which features were selected and their contributions to model goodness of fit as measured by the improvement in the coefficient of determination. After each model was constructed, we tabulated the subgroup and spatial structure of each explanatory variable selected to assess whether the variables selected within these analyses (Section 2.2-2.3) also improved the predictive power of the variable choices selected within the SVMR models. The first set of SVMRs used all variables other than dependent variable as possible inputs, with the goal of identifying relationships between individual variables. The second set used PC1 and PC2 from each disciplinary subgroup as possible inputs with the goal of identifying more generalizable flows of information from the major axes of variation within and between subgroups. In all cases SVMRs are used to identify

directional relationships between all possible pairs of variables (i.e., finding variable A is informed by variable B does not require B is informed by A).

Finally, we compared performance of the SVMRs selecting features from the full variable set to those selecting from a random subset. We constructed 100 SVMRs using 10 randomly selected features as possible inputs for each variable. We used one-way ANOVA and Kruskal-Wallis tests as a basis to assess performance differences between models with the full feature set vs. random subset, reporting  $p_{ANOVA}$  and  $p_{KW}$ , respectively. We interpret SVMRs selecting from the full feature set performing significantly better than those selecting from a random subset of features as confirmation that the methods are identifying relationships that are at least mathematically non-random.

## **2.5 Literature analysis**

To assess the presence and relative frequency of studies jointly considering relationships between each pair of variables in our data set, we conducted a series of searches using the Scopus database in October 2020, following methods from similar studies (Ward, 2015; Yoder et al., 2020). Each variable in our data set was assigned one or more keywords that are commonly used to describe that variable in the literature (Ward, 2021). Literature was searched for every pairwise combination of variables (12,246 unique searches) for studies containing both keywords and a required term to indicate a study was likely relevant to our study of river corridors (one of: river, stream, water, aquatic). We tabulated the total number of studies returned from each search to assess the interactions between variables that have been studied jointly with greater or lower frequency, and compared these results to the interactions found to be significant within the SVMR analysis. Conversely, we also assessed if the specific pairwise interactions identified as significant in the SVMRs were present in the literature.

## **3. Results**

### **3.1 Principal component analysis**

#### **3.1.1 Principal component analysis on all variables**

The PCA on all variables identified major axes of co-variation without regard to disciplinary grouping. PC1 explained 20% of the total variance (Table 2A), and contained mainly variables from the metabolomics subgroup, generally representing a gradient moving from terrestrially-derived aromatic compounds that are more thermodynamically favorable for microbial respiration to more microbially-derived compounds that are less thermodynamically favorable. PC2 explained 17% of the total variance and contained variables from the geologic setting subgroup, such as valley width and stream slope, showing marked gradients from headwaters to downstream reaches. Taken together PC1 and PC2 suggest that sampling sites within the river network are organized by organic matter chemistry and geology, which are themselves linked by terrestrial vegetation and soils.

### ***3.1.2 Principal component analysis on disciplinary subgroups***

PCAs were conducted on each subgroup to identify major axes of variation within individual disciplinary perspectives. The first two PCs within each subgroup explain an average of 52% of the within-group variance (median 46%, range 33-76%; Fig. 1A; Table 1). For physical chemistry, we interpret PC1 as representing weathering rate (from high to low) and PC2 as representing age of water (from high to low). For the geophysical setting, we interpret PC1 as representing network position (from headwaters to larger rivers) and PC2 as representing surficial geology. For nutrients, we interpret PC1 as representing enzymatic activity (low to high) which is itself the inverse of dissolved inorganic nutrient availability, and PC2 represents the accumulated organic matter in the shallow streambed. For metabolomics, we interpret PC1 as reflecting gradients from terrestrially-derived aromatic compounds that are more thermodynamically favorable for microbial respiration to more microbially-derived compounds that are less thermodynamically favorable. The metabolomics PC2 is interpreted as a gradient being dominated by products from organic matter degradation at one end and less-processed terrestrially-derived organic matter at the other end. For bulk DOM, we interpret PC1 as representing DOM quality from less to more humic or terrestrial in origin, and PC2 as representing microbial and proteinaceous DOM (from more to less). For macroinvertebrates, we interpret PC1 as representing richness (high to low) and PC2 as representing abundance (high to low). For stream solute tracers, we interpret PC1 as representing short-term storage of tracers

(low to high) and PC2 as representing the importance of advection and longitudinal dispersion to tracer transport (low to high).

**Table 1.** Result of principal components analyses conducted on all variables in a single analysis (top) and on each expert subgroup (bottom).

<b>PCA on all variables</b>						
	<b>PC1</b>			<b>PC2</b>		
	Variance explained (%)	Positive loadings	Negative loading	Variance explained (%)	Positive loadings	Negative loading
All variables	20	Nominal oxidation state of Carbon, % tannin, % condensed hydrocarbons, Modified aromaticity index, % Lignin	Gibbs free energy, % lipids, double-bond equivalency minus Oxygen, % protein	17	stream valley width, stream order, alluvium, valley width, discharge upstream, discharge downstream, advection-dispersion: MAD and D, segment sinuosity	valley segment slope, stream segment slope
<b>PCA on subgroups</b>						
	<b>PC1</b>			<b>PC2</b>		
	Variance explained (%)	Positive loadings	Negative loading	Variance explained (%)	Positive loadings	Negative loading
Physical Chemistry (PCHEM)	40 *	—	Mg, Ca	26 *	18O, 2H	—
Geologic Setting (GEO)	17 *	stream order, channel width, channel depth, segment sinuosity, alluvium, segment valley width, cobbly-sandy-loam	segment stream slope, segment valley slope, valley slope, stream slope	16	soil depth < 3 ft, % clastic flows, gravelly-clay-loam, greenish breccia residuum/colluvium, soil erosion severity, poor water yield	travel time to outlet, glacial drift, soil gravelly sandy loam, % soil depth 3-to-10ft, % ridge-capping lava flow, moderate water yield, live biomass
Nutrients and enzymatic activity (NUTS)	29 *	beta-D-glucosidase (C-acquiring), Leucine aminopeptidase (N-acquiring)	—	14	% Organic Matter in sediment	—
Metabolomics (ICR)	48	Nominal oxidation state of carbon, % tannin, % Condensed Hydrocarbons, Modified Aromaticity Index, % Lignin	Gibbs free energy, % lipids, Double bond equivalency minus Oxygen, % protein	28	% AminoSugars, % Carbohydrates	Aromaticity index, Double-bond equivalence
Dissolved Organic Matter (DOM)	47	peak A (humic-like), peak C (humic-like), total fluorescence	—	20	peak T (protein-like)	fluorescence index
Macroinvertebrates (MACRO)	30	—	Richness, Shannon, index, Richness of collector-gatherers, Richness of predators short term storage	16	Abundance of collector-gatherers	Abundance of shredders, Abundance of small body size
Stream Solute Tracer (TRACER)	19 *	—	(holdback, skewness, CV)	16	Dispersion, Fraction of mass in A/D, velocity, upstream and downstream discharge	—

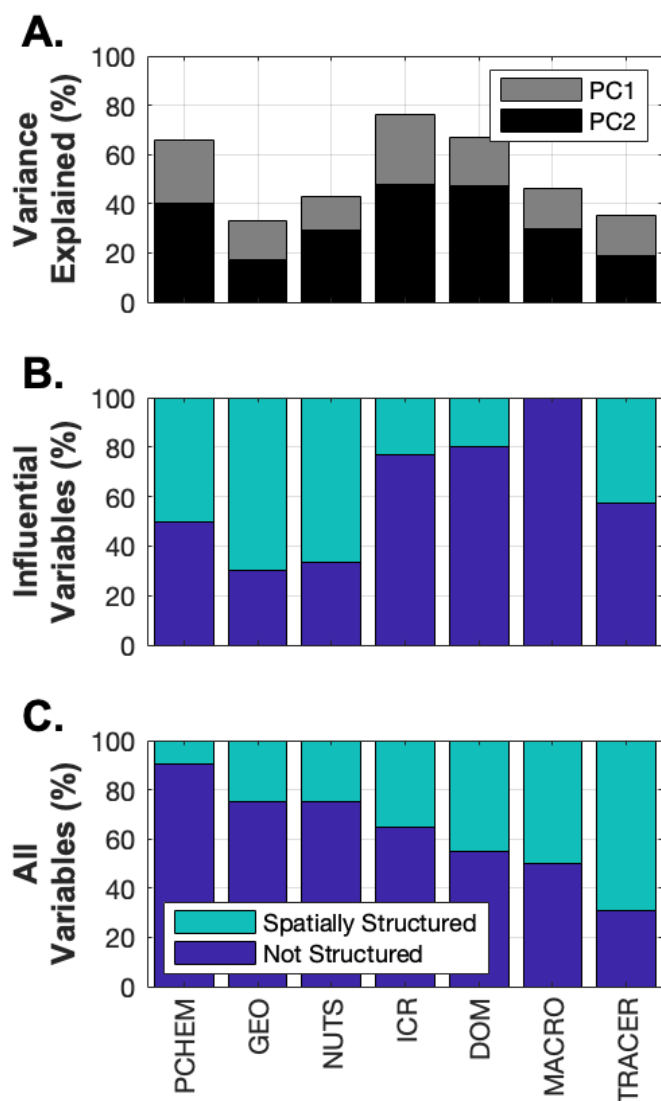
\* Indicates the PC is spatially structured

### 3.2 Spatial structure

Next, we assessed the degree to which variance in each variable can be explained by spatial structure. Of the 157 variables considered, we identified 56 variables (about 36%) as having spatial structure, compared to 101 variables (about 64%) without spatial structure. All structured variables were identified based on a linear semivariogram, with none exhibiting a spatial scale at which variation stopped increasing with distance between sample locations. This indicates variance in these spatially structured variables either (a) increases without bound or (b) only plateaus at scales that are larger than were included in the 5<sup>th</sup> order river basin we studied. This is consistent with prior studies of rivers, which exhibit fractality over a wide range of scales (e.g., Rodríguez-Iturbe & Rinaldo, 1997), with constraints (i.e., scale cutoffs) only occurring at relatively large scales (e.g., lateral valley constraints) and which may be functionally

unconstrained in the longitudinal dimension until they reach the ocean. Still others have found spatial structure in some parameters (e.g., in-stream solute concentrations) at scales that were encapsulated within our study (e.g., McGuire et al., 2014), suggesting that finding of spatial correlation lengths in one system or for one variable may not be universally transferable.

The fraction of influential variables with spatial structure was varied between subgroups (Fig. 1B, 1C), with 6 of 14 subgroup of PCs containing both structured and unstructured variables. The largest proportion of spatially structured variables were in the TRACER subgroup (69%; Fig. 1C), and the least were in the PCHEM subgroup (9.5%; Fig. 1C). The variables that appear in the disciplinary subgroup PCs did not separate into distinct groups of structured vs. unstructured variables. Instead, we found 44% of all influential variables were spatially structured (23% in PC1 and 21% in PC2) compared to overall 36% of all variables exhibiting spatial structure. All subgroups contained some structured influential variables except for MACRO (Fig. 1B), where only unstructured variables were selected.



**Fig. 1.** (A) Variance in the Andrews river corridor data set explained by PC1 and PC2 for each expert subgroup. (B) Percentage of influential variables (i.e., the variables included in the first two PCs) that do and do not have spatial structure. (C) Percentage of all variables within each subgroup that do and do not have spatial structure.

### 3.3 Support Vector Machine Regression (SVMR)

#### 3.3.1 Prediction of each variable using all other variables

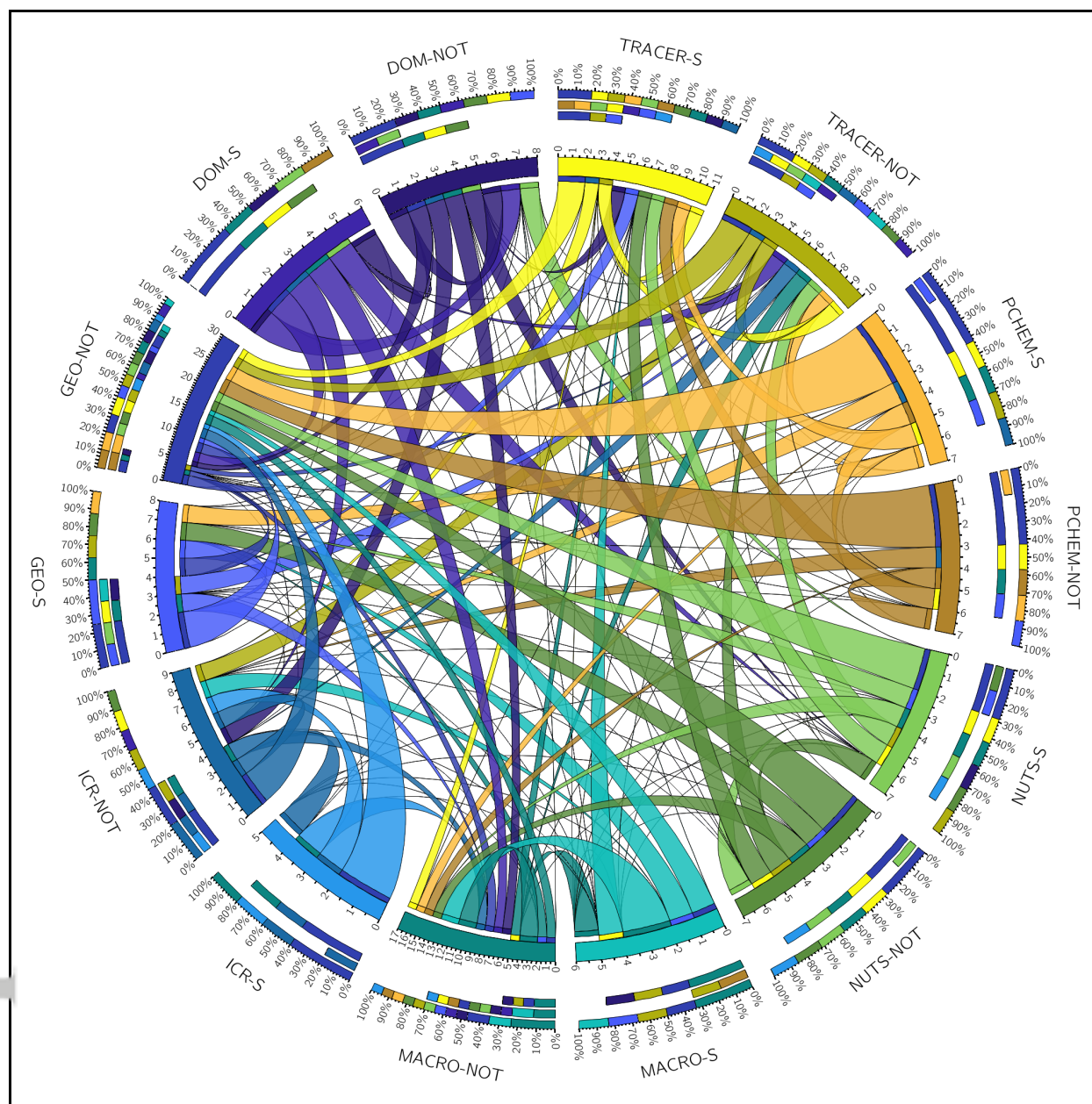
We identified 672 apparent relationships in the SVMR analysis that, taken together, demonstrate a complex network of interactions among variables in the river network, including variables that are typically measured by different research communities, and, hence, are commonly not measured at the same location (Fig. 2). The SVMRs were able to explain much of the variance

Accepted Article

in the underlying data, with an overall mean  $r^2$  of 0.83 (median 0.94, range 0.00 - 1.00). SVMRs for individual variables selected an average of 4.4 variables as predictors (median 4, range 1 to 10; Fig. S1), indicating that the relationships (i.e., statistical models) identified by the SVMRs were reasonably parsimonious. Additionally, performance of the SVMRs built from the full feature set was significantly better than those built from a random selection of features ( $p_{ANOVA} = 1E-19$ ;  $p_{KW} = 4E-29$ ), indicating SVMRs are selecting meaningful features and the associated relationships are appropriate for further analysis. The models built for spatially structured variables had an overall mean  $r^2$  of 0.91 (median 0.97, range 0.08 - 1.00) compared to a mean  $r^2$  of 0.78 for unstructured variables (median 0.90, range 0.00 - 1.00). Goodness of fit was also statistically better for the spatially structured variables ( $p = 0.008$ ; one-way ANOVA), indicating that spatially structured variables were more accurately predicted (i.e., higher  $r^2$ ) compared to unstructured variables.

Of the 157 variables predicted, 22% (34 variables) are informed by only out-of-group variables (i.e., variables from a different subgroup), and 11% (17 variables) are informed by only within-group variables (i.e., variables in the same subgroup). Thus, 67% of variables (106 out of 157) required both in-group and out-of-group information for optimal prediction by the SVMRs. Moreover, out-of-group information dominates predictor selection, representing an average of 59% of variables selected (median 66%, range 0-100%; Fig. 2, Table S1). Spatially structured variables represent an average of 27.3% of variables selected for individual SVMRs (Fig. S2A, S2C). Across the 157 SVMRs constructed, 30% (47 variables) did not select any spatially structured features. We found 3% of models (5 variables) selected only spatially structured features, and the remaining 67% (105 variables) selected a combination of structured and unstructured variables.





**Fig. 2.** Information flow within and among subgroups of variables commonly used as measures of river corridor dynamics based on the suite of SVMs constructed for each variable (Section 3.3.1). The variables included in the 7 subgroups (PCHEM = physical chemistry; GEO = geologic setting; NUTS = nutrients; ICR = metabolomics; DOM = dissolved organic matter; MACRO = macroinvertebrate; TRACER = stream solute tracer; variables in each grouping are detailed in Ward (2021)) are further organized by those with spatial structure (“-S”) and without spatial structure (“-NOT”).

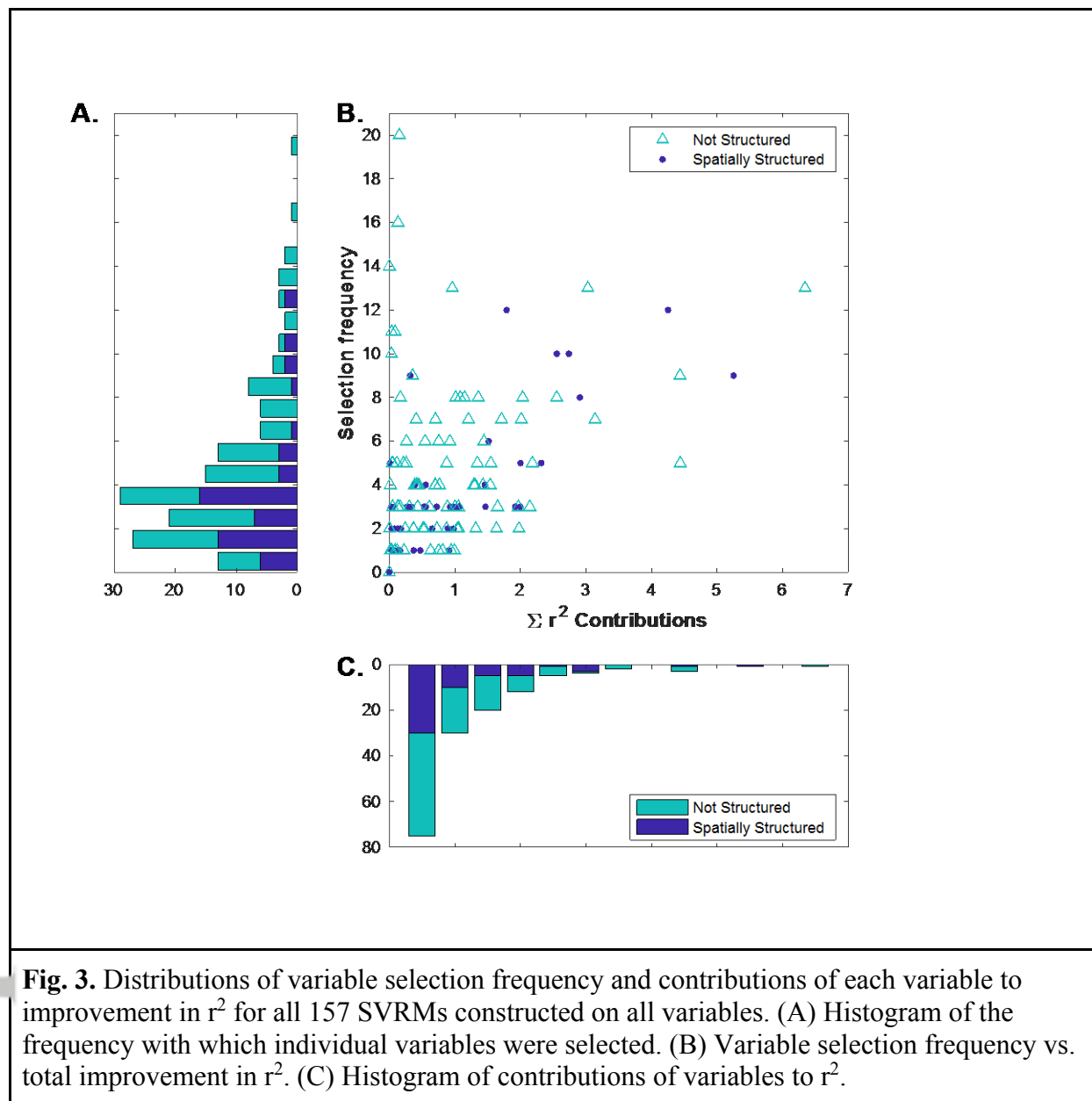
Each subgroup is represented by a different color to enable visualization of interactions with other subgroups, with the color of each ‘ribbon’ denoting the origin of information (i.e., the subgroup from which information flows). The width of each ‘ribbon’ denotes the relative frequency of interaction between variable groups.

The three ‘rings’ around the outside of the plot represent information flow between subgroups as:

- Inner Ring: destination(s) of information from each subgroup (i.e., answers the question “which other subgroups used information from this subgroup?”; colloquially the ‘outflows’ of information from one subgroup to another). These are the independent variables requires as inputs to make predictions of dependent variables in other groups.
- Middle Ring: the source(s) of information to a subgroup (i.e., answers the question “which variable informed relationship using to predict variables in a given subgroup?”; colloquially the ‘inflows’ of information to a subgroup). These are the independent variables providing information for predictions of variables within this group.
- Outer Ring: Scaled, total interactions with other variable groups regardless of directionality (i.e., answers the question “how related is this subgroup to others in the web of relationships?”). These are the relative magnitudes of direction-independent relationships between subgroups.

Individual variables were selected an average of 4.3 times (median 3, range 0-26; Fig. 3A). The most selected variable was in-stream  $\text{NH}_3$  concentration. However, this variable only contributed 0.046 improvement in  $r^2$  summed across the 26 models where it was selected. In contrast, the largest improvements were associated with the functional richness index for macroinvertebrate communities, which provided a total improvement of 6.3 in  $r^2$  summed across the 20 models where it was selected (average improvement of 0.315 in  $r^2$  when this variable was included in a model). Overall improvement associated with adding any variable was 0.83 (median 0.47, range -0.04 to 6.3; Fig. 3C).

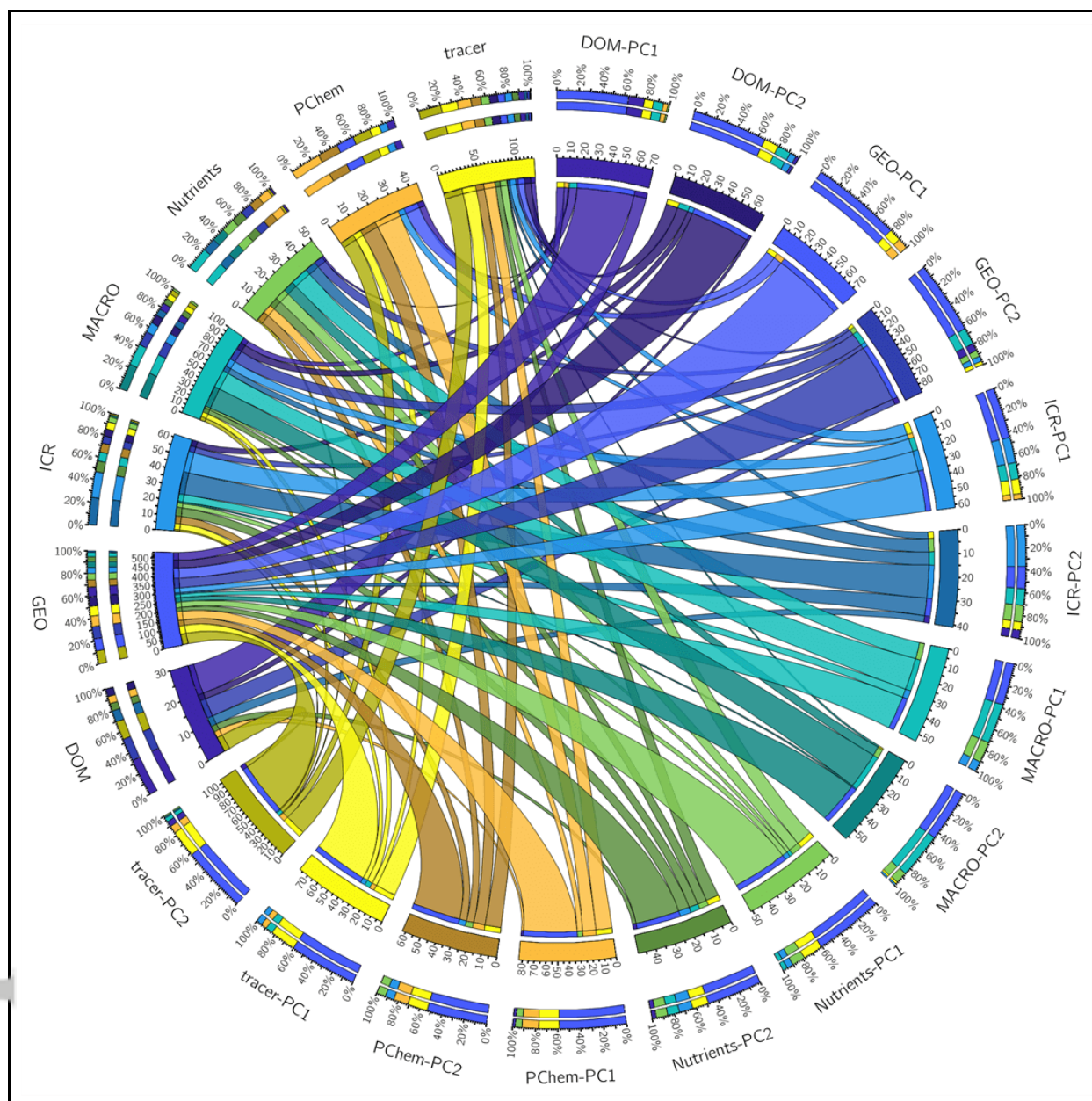
Across all 157 SVMRs constructed with the entire variable set, out-of-group variables were selected more frequently than within-group variables and contributed more to the overall  $r^2$  of the model. We found out-of-group variables represent about 30% of all selections within the SVMRs (Fig. S2C), but contribute more than 50% of the improvements in model performance (Fig. S2D). Similarly, spatially structured variables represent about 36% of all variables selected (Fig. S3C) and contribute about 40% of the improvements in model performance (Fig. S3D). These results indicate that river corridor variables typically considered to be outside the primary domain of individual field studies have a disproportionately larger effect than variables considered to be within the primary domain.



### 3.3.2 Prediction of each variable using principal components from each subgroup

The first two PCs for each subgroup define major attributes of the river network, as described previously in Section 3.1, but still leave an average of 48% of variance unexplained within each subgroup. To relate major axes of variation between subgroups, we constructed SVRMs for each variable using the PCs from each subgroup as inputs (Fig. 4). In-group PCs were always selected more frequently than PCs from any other subgroup (Table S2). In fact, about 25% of variables (39 of 157) were predicted solely from their in-group PCs. The explanatory power of PCs for in-

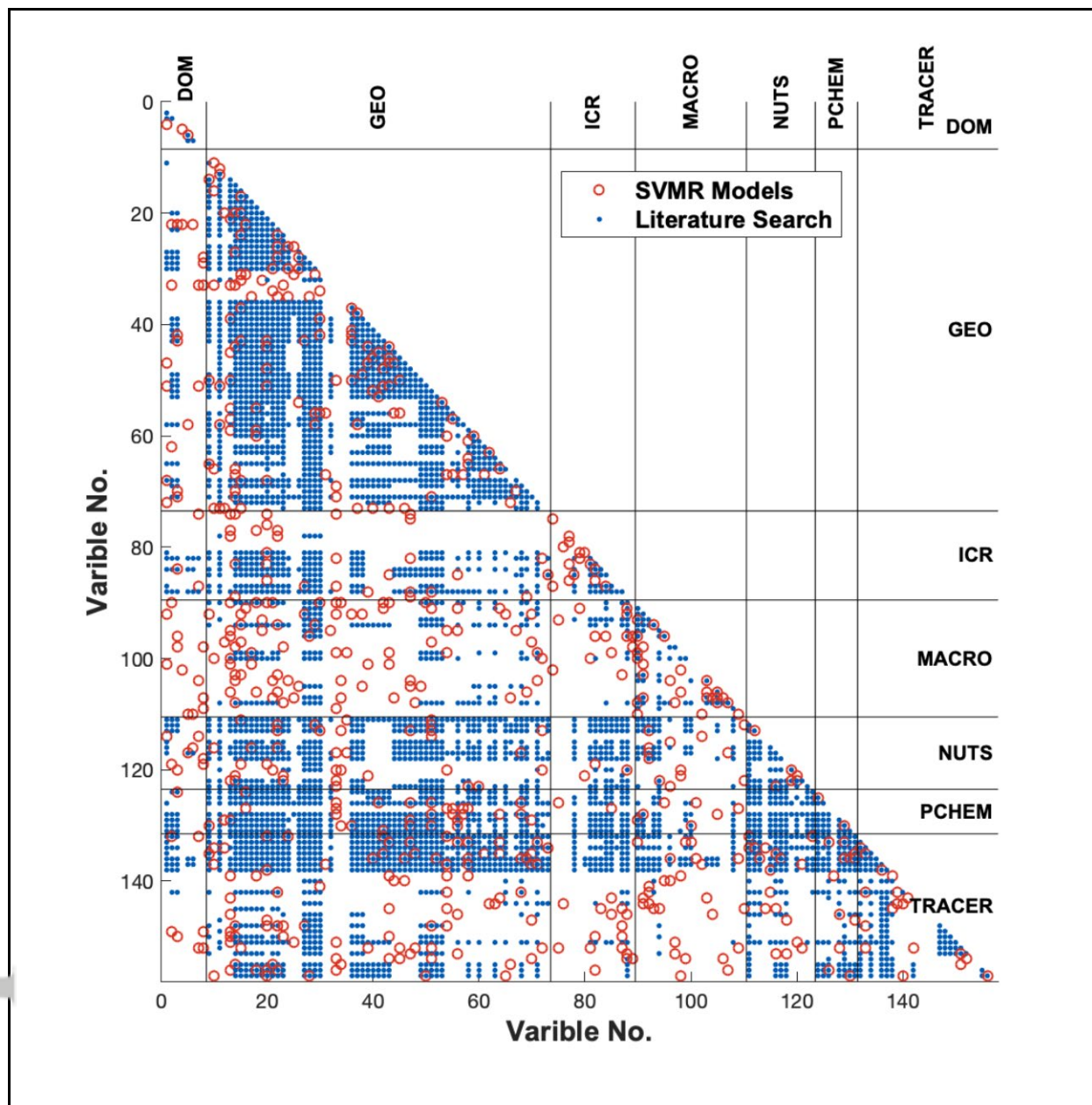
group variance is unsurprising given that PC1 and PC2 were successful in explaining an average of 52% of variance within their group. However, we also found about 26% of variable predictions (41 of 157) used only out-of-group PCs, and 118 variable predictions selected at least one out-of-group PCs. Further, variables in each subgroup drew information from nearly every other subgroup (see Table S1), These findings indicate that studies that are limited to one discipline are unlikely to explain as much of the observed variance in the measured variables as studies that intentionally span disciplinary boundaries, and that it is important for disciplinary understanding to at least characterize the major attributes from other subgroups.



**Fig. 4.** Circos plot showing the one-way flow of information from the subgroup PCs (Table 1; labeled “XXX-PCY” where XXX is the subgroup and Y in the PC number) to variables predicted by the suite of SVMRs described in Section 3.3.2. Plot layout and interpretation is identical to that described for Fig. 2, except that ‘flows’ of information only originate the PCs (i.e., subgroup PCs have only outflowing and total interactions; middle and outer rings) and only inform variables in the subgroups (i.e., variable subgroups only have inflowing and total interactions; inner and outer rings).

### 3.4 Studies of inter-relationships between steam corridor variables reported in the literature

Our literature search identified 4,075 combinations of variables that have been studied pairwise in the literature (of 12,246 possible combinations). The pairwise literature search returned a total of 2,731,694 results. The number of studies identified for any given pair of variables was highly skewed: 50% of published studies included the 18 most commonly studied pairs of variables (Ward, 2021), while the number of studies of any given pair of variables ranged from 1 to 270,015 (mean 670, median 14). These findings indicate a bias toward co-observation and reporting of a limited number of pairwise studies, consistent with a prior study that manually reviewed search results (Ward, 2015). We also found the existing literature is more focused on in-group relationships (57.2% of pairwise results) compared to between-group relationships (42.8% of pairwise results). In contrast, our SVMR approach identified a total of 672 pairwise relationships, of which 68.8% are between-group. Notably, about 84% or 564 variable pairs do not appear to have been reported previously (i.e., our systematic literature search did not return any manuscripts containing information on both variables). The remaining 16% (108 relationships) have been previously reported in the literature (Fig. 5). The 108 relationships found in both the literature and in our data analysis only represent about 2.6% of all previously-reported relationships, but these relationships are included in more than 16% of all published studies, indicating that prior studies have focused primarily on a relatively small number of relationships. On the basis of within- and between-group frequency, the literature is broadly not reflective of our findings, with the SVMR identifying higher frequencies of between-group relationships that are present in the literature (Table S3). Finally, we note that the lack of a relationship in the SVMR does not necessarily indicate that some relationship may be possible, just as the presence of a statistical relationship does not necessarily indicate a causal relationship. Some meaningful relationships could have been omitted due to signal-to-noise ratios, lagged correlation between variables, or because a highly correlated variable was already selected. This may explain why some well-studied relationships were not apparent in our analysis (Fig. 5).



**Fig. 5.** Scatterplot showing pairwise study in the literature (blue dots) and identification of a relationship in our SVMR approach (red circles) for all variable pairs. Variable numbers correspond to the order variables are listed in Table S1.

## 4. Discussion

### 4.1 Relating large-scale spatial patterns and localized heterogeneity in the river corridor

Spatial structure alone is not sufficient to explain the inter-relationships between variables that we observed in the river corridor. We found that spatially structured variables were included in SVMRs less frequently than would be expected by random chance (i.e., structures variables are

27% of the variables included by SVMRs although they make up 36% of the total variable set). This means the predictions of spatially structured variables were not dominated by structure from a small number of structured variables. Further, a majority of variables observed (about 64%) were not themselves spatially structured, and five of the seven subgroups (PCHEM, GEO, NUTS, ICR, TRACER) resulted in at least one PC that was not spatially structured. These results indicate that spatial structure is not ubiquitous in the river corridor. Instead, some variables represent local ‘noise’ on the network-scale ‘signal’ (i.e., systematic variation in physical, chemical, and biological processes from headwaters to large rivers; Vannote et al. 1980). This heterogeneity is either independent from large-scale system structure (i.e., controlled by local process interactions that are neither controlled by nor influence larger scale patterns) or simply have sufficiently high variability to obscure larger-scale trends. Such localized ‘noise’ may also reflect processes whose importance is localized in space or time, but do not recognizably follow a larger spatial structure.

Individual variables reflect complex interactions that can either lead to the emergence of spatial structure or overwhelm the underlying spatial structure associated with more basic variables like slope and elevation. We found six variables that were spatially structured but had strong relationships (SVMRs) that only included unstructured variables. In these cases, spatial structure emerged or was generated by the interaction of variables that did not themselves have spatial structure. Conversely, 60 of the SVMRs for unstructured variables included at least one spatially structured variable (38 selected 1, 14 selected 2, and 8 selected 3 spatially structured variables). This pattern suggests that spatial structure does not necessarily propagate from one variable to another, indicating “signal shredding” in the river corridor (Jerolmack & Paola, 2010), where information is erased by interactions between variables. While such behavior has only been confirmed previously for sediment transport, our findings indicate that localized feedbacks can generally overwhelm underlying spatial structure within the river corridor. This suggests that sufficiently large perturbations will have system-wide impacts (e.g., large fires, floods), but internal dynamics may overwhelm large-scale patterns under normal circumstances. Consequently, studies of river corridors must consider local-scale interactions (i.e., internal dynamics), large-scale drivers (i.e., external forcing), and the temporal context (i.e., historical contingencies) if we are to account for the feedbacks and interactions in the river corridor.



## ***4.2 Benchmarking inductive relationships to established, deductive science***

While a majority of the apparent relationships identified in the SVMR are novel compared to the literature, the inductive approach did identify a suite of relationships that are consistent with pre-existing conceptual models from the literature and published findings from the H.J. Andrews Experimental Forest. Below we detail three examples of consistency between inductive and deductive science in the basin, including relationships that are generally viewed as important in the river corridor: hydrologic exchange processes, microbial ecology, and the River Continuum Concept (Vannote et al., 1980). Taken together, these examples demonstrate that our inductive approach is able to extract meaningful relationships from data, building confidence that never-before-reported apparent relationships are worthy of future study. The inductive identification of patterns and couplings that are consistent with deductive work, and presented in subsequent subsections, is important as it confirms that meaningful relationships can be extracted from complex data using inductive approaches.

### ***4.2.1 River Corridor Exchange***

Our findings indicate that geologic setting, and the resultant land cover and soils, are important controls on solute transport patterns in the river network. In prior analysis, we focused on spatial patterns in reach-scale solute transport and identified substantial, unexplained heterogeneity in univariate regressions (Ward, Wondzell, et al., 2019). The SVMRs in this study included 35 unique variables that predict the 11 observations that common to our analysis and the prior work. These variables primarily fall within the geologic setting ( $n = 10$ ), tracer (8), and macroinvertebrate (7) groups. Of those variables, the abundance of the oldest exposed lava flows was included most commonly (5), followed by slope stability and forest cover (3 each). Five additional variables were selected twice (two associated with geological setting, two with tracer, and one with macroinvertebrates), while 26 variables were selected by only one SVMR. Notably, geologic setting was selected more frequently than other descriptors of tracer transport, suggesting autocorrelation amongst metrics describing tracers is not sufficiently strong to overcome the heterogeneity imparted by the landscape. This finding is in good agreement with several prior studies that have identified geologic setting as a high-level control of river-groundwater interactions and hydrologic travel time based on results from both field

observations (Payn et al., 2009; Valett et al., 1996) and models (Cardenas, 2008; Frissell et al., 1986; Wondzell & Gooseff, 2014; Wörman et al., 2007).

Ward et al.'s (2019) observation of monotonic trends between most hydrologic exchange metrics and discharge - which they describe as a proxy for network position - agree with our finding of spatial structure in several variables describing geomorphic setting (including hydraulic conductivity, valley slope, valley width, sinuosity), river flow (velocity, discharge), and solute transport metrics (e.g., median travel time, skewness). We did not find spatial structure for other metrics of exchange where Ward et al. did, including the coefficient of variation, holdback, and channel water balance. Further, many of the relationships identified by Ward et al. have low explanatory power as evidenced by low  $r^2$  values, indicating that hydrologic exchange cannot be described by a single explanatory variable. In contrast, the multivariate and nonlinear responses encoded in the SVMs better explain the patterns in river corridor exchange observed in the Andrews watersheds.

#### **4.2.2 Microbial Community Assembly**

Interactions along the river corridor can not only 'shred' or erase information (*sensu* Jerolmack & Paola, 2010), but can also generate new information and patterns. For example, prior work at the H.J. Andrews Experimental Forest spanning headwaters through 5<sup>th</sup> order rivers (Wisnoski and Lennon, 2021) showed that microbial assemblages in headwater streams habitat-dependent, while the microbial community became more homogeneous with distance downstream.

Additionally, the same study found taxonomic  $\beta$ -diversity was explained by an axis with positive loadings for elevation and dissolved organic carbon, and negative loadings for electrical conductivity, pH, total nitrogen, and total phosphorus. Microbial assemblages are known to arise in response to local heterogeneity in the landscape, integrating inputs and environmental variables in space and time. While we did not analyze microbial assemblages explicitly here, we can interpret our observations in the context of prior findings at the site (Wisnoski and Lennon, 2021). Our results show spatial structure in electrical conductivity and several geomorphic variables that are known to vary with elevation, but no spatial structure in total dissolved phosphorus, DOC, or total dissolved nitrogen. In comparison to the controls on taxonomic  $\beta$ -diversity described by Wisnoski and Lennon (2021), we did find spatial structure in elevation, in-

stream nitrate+nitrate, and electrical conductivity, but not in bulk dissolved organic carbon, ammonia, or total phosphorous. Thus, our findings are broadly consistent with past findings that at least some of the known controls on microbial diversity are spatially structured. However, we also note that not all controls were structured, but the related microbial community did retain spatial organization. Consequently, studies focused at single locations along a stream may be missing contextual information on controlling factors that have propagated from the catchment headwaters, or misinterpreting signals that were generated within the river corridor itself.

#### **4.2.3 River Continuum Concept**

The River Continuum Concept (Vannote et al., 1980) -- one of the most widely recognized and cited conceptual model of river corridors -- argues that Leopold's conceptual model that geomorphology reflects energy equilibrium can be extended into ecosystem functions (Langbein & Leopold, 1966; L B Leopold et al., 1964; Luna B. Leopold & Langbein, 1962). Vannote et al. (1980) specifically proposed: (a) biological communities should evolve to optimize the use of available energy (i.e., biodegradable organic matter); and (b) energy availability will vary systematically from headwaters to large downstream rivers. Our PCA results on all variables are broadly consistent with these hypotheses, which is to be expected at the H.J. Andrews Experimental Forest was one of the key sites studied in developing and demonstrating the conceptual model. We found organic matter chemistry and geological setting explained 37% of the variance across the entire data set (PC1 and PC2; Table 1). With regard to biological communities optimizing to use available energy in an organized fashion, we do see that available energy – in this case assessed via energy availability in organic carbon (PC1 on all variables) – defines one critical dimension of variation in the system. Additionally, the high proportion of spatially organized variables in TRACER, GEO, and NUTS is consistent with broad concepts of systematic organization along river networks. Indeed, we found spatial structure in about 36% of all variables across all disciplinary subgroups, consistent with the idea that large-scale gradients will drive systematic trends in both physical and biogeochemical processes. We did find spatial organization in shredders which is consistent with the River Continuum Concept. Our findings on the importance of organic carbon as an explanatory variable for patterns in the river corridor also support Vannote et al.'s expectation of the importance of energy availability to the structure of fluvial ecosystems.

### ***4.3 Open questions stemming from the inductive analysis***

We applied machine learning techniques to cross-disciplinary data to uncover novel relationships that are worthy of subsequent investigation. Inductive approaches cannot reveal causal relationships, making this a useful approach to identify relationships for future study, rather than proving mechanistic pathways. To demonstrate the value of this approach, we explore a selection of findings from the network of relationships identified by our SVMR models, focusing on relationships that are at the cutting edge of our understanding of river corridors. While our body of knowledge has methodically built knowledge and is beginning to engage with these questions, we take it as a positive sign that inductive approaches were able to also pick these relationships out of the data set. Thus, in addition to consistency with past findings (Section 4.2) we take these findings as further support that inductive approaches are able to identify relationships worthy of further scrutiny. We pose these as potential areas for future study to highlight the role of inductive analysis as a path to inspire the asking of questions, rather than providing mechanistic answers, about the complex structure and function of river corridors.

#### ***4.3.1 Why are metabolomics data most informed by geological variation?***

Metabolomics data alone formed PC1 for the overall analysis, explaining 20% of the variation in all data analyzed (Table 1), while geomorphic variables dominate PC2, explaining 17% of all variance. Moreover, these axes are, by definition, orthogonal implying that the two groupings should be independent. Across the 16 SVMRs constructed on organic carbon chemistry (ICR) variables, none selected any features from the dissolved organic matter, nutrient, nor physical chemistry subgroups (DOM, NUTS, and PCHEM, respectively). Instead, out-of-group information was exclusively from geological features, solute tracer, and macroinvertebrate groupings (GEO, TRACER, and MACRO, respectively). This is particularly surprising given that a host of variables traditionally used to describe organic matter were available, including optical measures of carbon quality (e.g., EEM features, SUVA<sub>254</sub>) and quantity (e.g., total DOC, carbon acquiring extracellular enzymes). We posit that the apparent dominance of physical setting over biogeochemical variables emerges through the microbial community (i.e., the Baas Beeking hypothesis; *sensu* O'Malley, 2008; Fondi et al., 2016; Wit and Bouvier, 2006). In other words, geologic setting and hydraulics set a template that defines which microbial communities

will occur, and these communities are responsible for the molecular form of organic matter that is transformed within and exported from a given location. This is, functionally, the River Continuum Concept applied to microbial communities. We expect the role of microbial community structure in defining ecosystem processes will be critical as we transition from conceptual models based on bulk measurement of organic matter (e.g., DOC, EEMs) to models informed by metabolomics.

Previously developed theories based on bulk DOC or proxies for organic matter quality must be revisited, because the field of metabolomics is rapidly evolving. The limited suite of studies that include both organic carbon chemistry and nutrient data (ICR and NUTS) make comparisons for consistency of findings limited. It is possible that previous conclusions about carbon limitations in some systems may have been biased by only considering bulk DOC or DIC instead of its molecular composition, which is highly nonuniform in its ecological function. We do not expect that organic matter molecular composition is entirely controlled by geologic setting (though such control has been reported; e.g., Robertson et al., 2019; Cotrufo et al., 2013), but instead that in-stream organic matter reflects the integration of physical, chemical, and biological processes occurring upstream of the sampling location. These processes are diverse, spanning the influences of terrestrial vegetation, soil-forming processes, photochemistry, organo-mineral interactions, and in-stream biological production and transformation of organic molecules. Thus, the core questions are to understand when, where, and how organic matter is produced, transformed, and transported. We expect that understanding microbial communities and their metabolism will be critical to answering these questions.

In addition, Danczak et al. (2020) proposed a conceptual framework that draws parallels between organismal birth, death, and dispersal and organic matter production, transformation, and transport. They argue that organic molecules are assembled into metabolomes via a combination of production, transformation, and transport just as organisms are assembled into communities via a combination of birth, death, and dispersal. Danczak et al. (2020) also provide an analytical approach for quantifying assembly processes, including the ability to infer when transport overwhelms influences of production and transformation. This approach may be fruitful in linking upland dynamics to aquatic dynamics (Waring et al., 2020; Wisnoski et al., 2021),

linking microbial community assembly processes to organic matter assembly processes, and further highlights the need for conceptual synthesis in the river corridor (Stegen et al., 2018).

Finally, metabolomics data has been used previously to inductively reveal limitations of using bulk water chemistry in river corridors to understand specific biogeochemical conditions. For example, there has been a recent revelation that conceptual models for denitrification in river corridors were framed at a large river network scale and not capturing dynamic, small scale controls of anaerobic metabolic pathways, including denitrification (e.g., Briggs et al., 2015). Since this revelation, field experiments and deductive methods have revealed that denitrification is in fact occurring in sediment “microzones” across a wide range of river corridor conditions that was previously hidden by and assumed impossible based upon bulk water chemistry (e.g., Knapp et al., 2017; Hampton et al., 2019; Hampton et al., 2020).

#### ***4.3.2 What controls nitrogen-acquiring extracellular enzymatic activity in a nitrogen-limited ecosystem?***

Aquatic ecosystems at the H.J. Andrews have been historically considered to be nitrogen limited (Sollins et al., 1981; Triska et al., 1984). Consequently, we expected that microbes would generate both leucine aminopeptidase (LAP) and N- acetylglucosaminidase (NAG) to acquire nitrogen and that this would be ubiquitous across the basin. Moreover, C:N:P ratios of extracellular enzymatic activity (EEA) should indicate an overproduction of N-acquiring enzymes as N-limited microbes allocate energy to acquiring their limiting nutrient (e.g., Sinsabaugh et al., 1997) .

To test this expectation, we considered two nitrogen-acquiring enzymes: LAP and NAG. LAP was part of PC1 for the NUTS subgroup and was orthogonal to total organic matter in the sediment, indicating little control on sediment organic matter in explaining LAP. SVMRs for LAP identify several GEO variables (bedrock type, hillslope stability, and channel water balance), allochthonous inputs to the river (deciduous forest, abundance of collector-gatherer macroinvertebrates), and organic carbon (spectral slope and ICR ‘other molecules’). Positive correlations with spectral slope and small molecules in the ICR indicate increased LAP occurs where relatively small and non-aromatic carbon sources are present. Similarly, NAG was

Accepted Article

predicted by bedrock type, ICR (protein abundance), and phosphorus-acquiring enzymes. Because we do not see spatial structure in LAP, NAG, nor 11 of the 13 variables selected by their SVMRs, we infer that there is not a spatial control on nitrogen acquiring enzymes.

Several studies have reported increasing EEA with nutrient availability (Hill et al., 2010; Sinsabaugh et al. 1997; Williams et al. 2010; Williams et al. 2012), which is not consistent with our findings (i.e., no measurement of bulk nitrogen, carbon, phosphorus, nor oxygen were selected by SVMRs for the ICR subgroup). Instead, we find that EEA may be explained by particular classes of organic matter – specifically smaller, less aromatic carbon molecules, consistent with Williams et al. (2012) and Hill et al. (2010). We also hypothesize the prevalence of GEO features selected by SVMRs but lack of spatial structure may indicate that there are geogenic micronutrient controls on the localized enzymatic activity that have not been measured, such as the availability of potassium, manganese, iron, and silica that weathers from local features.

Another enzymatic question that requires more deductive work is whether the entire river corridor is N-limited. Ecoenzymatic ratios of 1:1:1 C:N:P suggest an equilibrium between microbial biomass and detrital organic matter (Sinsabaugh et al., 2009). The ratios of C:N and C:P acquiring enzymes in our study (GLU:LAP+NAG and GLU:AP, respectively, based on data in Ward et al., 2019) have slopes that are statistically indistinguishable from analyses of global datasets (Sinsabaugh and Shah, 2012), indicating EEA is produced in relative proportions to the basic C:N:P ratios required by microbes, suggesting that the sediment microbial community may not, in fact, be N-limited relative to the availability of other nutrients and substrates. Therefore, while catchment-scale mass balances indicated one understanding of the system as N-limited (e.g., Sollins et al., 1981; Triska et al., 1984), we interpret the EEA data as an indicator that the microbial community has adapted to the available N, and that this is present across the network (based on the lack of spatial structure).

Our analyses suggest many fruitful paths forward for interdisciplinary river corridor research. These include, but are not limited to, the examples presented above that (a) relate molecular characterization of carbon to EEA to investigate organic matter quality controls; (b) comprehensively sample stream, streambed sediment, hyporheic pore water, and hyporheic

Accepted Article

sediment communities for EEA to test our hypotheses that microbes are not N limited across these spatial domains; and (c) use repeated measurements to assess if one spatial snapshot of the network adequately captures temporally dynamic behavior (as was found in Giraldo et al., 2014). Our findings also suggest that the concept of ecological stoichiometry and nutrient limitations manifest differently across multiple scales, warranting consideration of the places, times, and scales at which equilibrium or limitation should be inferred, and whether findings of limitations at one scale can be directly transferred to other scales. One particularly compelling question resulting from our work is whether system-wide, large-scale N-limitation indicate low N inputs at all scales, internal limitations due to spatial structure or heterogeneity (e.g., localized inputs from N-fixing alders), biogeochemical limitations (e.g., kinetics of organic matter breakdown), or transport limitation (e.g., inaccessibility of nutrients in some locations)?

#### **4.4 Inductive relationships are observations around which hypotheses can be spun and tested**

The suite of models we constructed include 672 apparent relationships, 84% of which have not been previously studied based on our literature search. It is important to recognize the relationships identified here are intended as future directions, not as endpoints that reflect a causal or mechanistic understanding, particularly in the case of correlations that have not been reported by other studies. Each relationship serves as a set of observations, the first step in the scientific method. We envision the next step for each relationship being the generation of hypotheses that propose mechanisms or explanations, followed by rigorous investigation with deductive approaches to rule out spurious correlation and other errors. While we have now used a coarse sieve to identify mathematically meaningful relationships in the data, additional study is needed to test the validity of each apparent relationship.

Even without additional investigation, it is perhaps surprising that so many apparent relationships identified by our inductive approach were not found in the literature search. Critically, without future study of hypotheses that can explain each relationship, like the few explored in Section 4.3, we cannot differentiate if the relationships are meaningful or spurious. In this regard, the inductive approach has fulfilled the promise of sieving nearly 25,000 potential relationships and identifying the 672 that warrant further scrutiny. While 108 of these have been



previously reported in the literature, we identify four possibilities to explain the lack of consideration of the remaining 564 pairwise statistically significant couplings in prior studies, and reflect on how these results can be used to advance our goal of synthetic science to yield comprehensive descriptions of the structure and function of river corridors.

#### ***4.4.1 Disciplinary, deductive science is the predominant mode of inquiry***

The norms of classical research funding opportunities and publications require deductive approaches, where the limited resources of time and financial support are focused on testing specific, mechanistic hypotheses. Consequently, researchers tend to dedicate effort and resources on a narrow suite of specific observations rather than broader datasets that may inform the connections between disciplines and scales. However, this paradigm is shifting with emphasis on macrosystems research (Heffernan et al., 2014), the explicit design of networks to facilitate synthesis (e.g., AmeriFlux, NEON, Critical Zone Collaborative Networks), and new funding initiatives. Our results show that the inherent complexity of river corridors and networks means that experimental programs of limited scope will often miss important process controls. This finding provides further support for our earlier recommendation that all river corridor studies collect a standard set of observations for fundamental system characterization (Ward, 2015), as this information is likely to be important to testing hypotheses in ways that may not be apparent in the initial study design. In this context, the inductive approach we propose here is extremely useful for rapidly identifying relationships spanning disciplinary boundaries that would otherwise take decades of disciplinary inquiry to identify.

#### ***4.4.2 Existing data sets are incomplete and could not have uncovered relationships***

Our analysis relies on the most comprehensive catchment-scale observations of interacting physical, chemical, and biological processes in any river corridor to-date. The dataset we analyzed also builds upon extensive prior work and data from the H.J. Andrews Experimental Forest. Such comprehensive datasets, particularly co-located with long term ecological research, have not previously been available and require extensive interdisciplinary collaboration to obtain. For example, molecular organic matter chemistry (e.g., FTIRCMS) is only recently emerging as part of river corridor science (Graham et al., 2018; Stegen, Johnson, et al., 2018; Zhou et al., 2019) and has not been jointly collected with the breadth of observations we

analyzed here. To make further progress in unraveling the complexity of river corridors, we recommend combining standardized system characterization across many streams and rivers with intensive study of select watersheds to generate the rich datasets needed to evaluate process interconnections and scale dependencies (Stegen & Goldman, 2018). In this case, the comprehensive nature of the data set explains why novel relationships were identified here: such breadth of data were simply not collected in past efforts. This further demonstrates the utility of inductive analysis in generating hypotheses from new datasets that can then be tested more broadly. Finally, note that our own data set, while comprehensive, is far from complete in terms of all variables that could be measured across all relevant spatial scales, temporal scales, and process dynamics.

#### ***4.4.3 Relationships may be scale- or time-dependent***

Both the structure and function of river corridors are known to be scale-dependent (Frissell et al., 1986; Rodríguez-Iturbe & Rinaldo, 1997; McCluney et al., 2014). The network scale considered here is larger than many studies of river corridors (see reviews by Tank et al., 2008; Ward, 2015). It is possible that the relationships identified between variables here by SVMR do not hold at all scales, or that the relationships are real but have not been tested over the range of scales we included in our analysis. Prior studies of river structure have found that self-similarities and scale dependencies generally only occur over a limited range of scales, and either average out at large scales or are limited by a physical constraint (e.g., water depth, channel width, valley width) (Jerolmack & Paola, 2010; Nikora & Hicks, 1997; Rodríguez-Iturbe & Rinaldo, 1997). As with relationships between individual variables, scale dependencies and scaling limits identified from broad data analysis must be considered as hypotheses and tested using directed observations and/or simulations with competing or alternative formulations. Similarly, analyses here focused on a data set collected under baseflow conditions and process controls are expected to vary in response to seasonal and storm dynamics in forcing. Moreover, our analysis are focused on what can be gleaned from a single snapshot in time, whereas the actual characterization includes a combination of variables spanning relatively dynamic (e.g., dissolved oxygen) to relatively static (e.g., valley slope), which may cause some relationships to manifest and obscure others. Future efforts to combine high temporal resolution data with spatial synoptic campaigns could directly address this limitation.

#### ***4.4.4 Spurious correlation may have driven the inductive relationships identified***

The relationships identified in our study may represent spurious correlation of disparate data or other mutual dependencies in the underlying data, a known limitation of machine learning approaches. In this case, the inductive approach aids in identifying mathematical artifacts rather than causal pathways or process interactions. Such relationships could also reflect redundant information (i.e., several different variables may reflect similar features on the landscape, and the autocorrelation amongst independently-measured variables may obscure underlying relationships). For example, if geology, land cover, and soils all systematically vary with increasing elevation, then these variables will all show consistent relationships that may confound interpretation. We emphasize here the relationships identified by SVMR and other machine learning methods only provide a starting point for generation of hypotheses, not an endpoint. The next step for investigation of such putative relationships would be to hypothesize a causal mechanism and design a study to collect the specific data needed to test it, while still capturing the essential system information identified here for purposes of evaluating scale dependency and complex system controls.

#### **4.5 Toward a unified conceptual framework for river corridors**

A unified conceptual framework for river corridors will require studies to move beyond the discipline-specific and site-specific studies that have dominated our field in the past decades (Ward, 2015; Ward and Packman, 2019). Instead, we need to augment our existing body of knowledge with ‘connective tissue’ that allows integration of our findings across spatial scales, temporal scales, and processes. Here, we endorse the conceptual organization Stegen et al. (2018) posed for microbial ecology, where we can begin to arrange our past and future studies around external forcing, internal dynamics, and historical context to explain and predict both temporal-variability and resultant services and functions of river corridors. Indeed, the framework of separating external forcing from internal dynamics is consistent with emerging theories in catchment hydrology where the same language has been applied to river corridors (Harman et al., 2016). However, this organization ultimately requires consideration of our studies in a synthetic framework rather than from a disciplinary framework.

Our study suggests that one avenue toward progress in river corridor science, complementary to common deductive approaches, is through the collection of uniform metadata and even observations typical of other scientific domains as part of disciplinary studies. We demonstrate here that, in the dataset we collected, out-of-group (i.e., cross-disciplinary) data were important to explaining many of the disciplinary (i.e., in-group) patterns that were observed. Thus, the out-of-group data not only enable synthesis, but also simultaneously improve disciplinary understanding by facilitating the generation and testing of new hypotheses. While the concepts of uniform metadata and common observations have been previously called for (Ward, 2015; Ward & Packman, 2019), our study demonstrates the value of these data to improve prediction of individual variables or functions in the river corridor. One potentially valuable path forward would be comprehensive characterization of several river corridors and at multiple times of year (i.e., a modern and disciplinary broader take on the work underpinning the River Continuum Concept; Minshall et al., 1983) to help determine which of the relationships we putatively identify here are fundamental and general, spurious, time-variable, or organized by larger climactic or geologic patterns. Another useful approach would be to identify and collect a small number of variables that are informative across many sub-disciplines, and organize the findings into spatially and temporally comprehensive datasets (e.g., Tiegs et al., 2019; Stegen and Goldman, 2018).

In this study, we have demonstrated an application of machine learning approaches to generate relationships that may inspire new studies to reveal the ‘connective tissue’ linking our understanding across spatiotemporal scales and disciplines. Indeed, the step of organizing raw observations to develop testable hypotheses is at the core of the scientific method, and we have prototyped one approach to organize observations and highlight potential relationships in the data. Hypothesis generation is touted as one of the core values of field-based observation and monitoring (Burt & McDonnell, 2015; Lovett et al., 2007), where observations demand explanations. The inductive approach used here presents a body of putative relationships for subsequent study, at least some of which are consistent with prior conceptualizations and observations of river corridors (section 4.2) and emerging areas of inquiry (section 4.3). We do not propose that such approaches supplant deductive science, but rather that the two approaches

must be coupled in river corridor science. The inductive approach provides an unbiased or naive data synthesis, which has the potential to reveal patterns and relationships that would not be obvious from our present, disciplinary perspectives.

## **5. Conclusions**

We began with the assumption that all variables may interact with all other variables, yielding nearly 25,000 potential pairwise relationships between variables. Using machine learning, we rejected most of these relationships, identifying 672 apparent relationships that have explanatory power in the data set, notably including 564 pairwise relationships that were not previously explored in the literature. Put another way, we have generated a web of 564 new apparent relationships that may reveal new couplings in the river corridor. These relationships eschew disciplinary or method-specific approaches, providing ‘connective tissue’ between traditional discipline-, scale-, site-, or method-dependent knowledge. Moreover, the network of relationships we have identified is consistent with several past studies from the field site (Vannote et al., 1980; Ward, Wondzell, et al., 2019; Wisnoski & Lennon, 2021), providing confidence that at least some of these relationships are more than spurious correlations.

Most of the relationships we identified, including a majority of those not present in the literature, include between-group flows of information. Our results show that interactions between processes that are typically studied by different disciplines is critically important to explain structure and function in the river corridor. This conclusion is, perhaps, unsurprising as a macrosystems view would acknowledge and expect to find cross-scale and interdisciplinary relationships (Heffernan et al., 2014; McCluney et al., 2014). Still, this view is seldom fully captured in existing experimental designs and the resulting data sets and literature. Importantly, we also demonstrated that spatial structure can be both generated through the interaction of unstructured data as well as destroyed or overprinted along the network. Thus, consideration of how an observed pattern may emerge or not be visible along a spatial gradient is a critically important consideration prior to interpretation of data sets.

Building connections between existing studies requires explicitly planning for synthesis in future efforts. Here, we demonstrated the value of collecting data sets that enabled synthesis within and

between locations, disciplines, and scales. This does not diminish the value of traditional, disciplinary hypothesis testing and deductive approaches to science. Instead, common metadata and even a small number of out-of-group observations may enable synthesis efforts based on inductive approaches that aids in spinning new hypotheses. Ultimately, inductive approaches are a useful way to generate hypotheses from existing observational datasets and advance our scientific understanding.

**Acknowledgements.**

This research has been supported by the Leverhulme Trust (Where rivers, groundwater and disciplines meet: a hyporheic research network), the UK Natural Environment Research Council (grant no. NE/L003872/1), the European Commission, H2020 Marie Skłodowska-Curie Actions (HiFreq, grant no. 734317), the U.S. Department of Energy (Pacific Northwest National Lab and DE-SC0019377), the National Science Foundation (grant nos. DEB-1440409, EAR-1652293, EAR-1417603, and EAR-1446328), the University of Birmingham (Institute of Advanced Studies), and with resources from the home institutions of the authors. Data and facilities were provided by the H. J. Andrews Experimental Forest and Long Term Ecological Research program, administered cooperatively by the USDA Forest Service Pacific Northwest Research Station, Oregon State University, and the Willamette National Forest. In lieu of detailed author contributions, we report that this study was conceptualized approximately 10 years ago and has benefited tremendously from discussions with a broad group of friends and collaborators. Work on this manuscript was initiated at the slow freshwater science meeting hold in Santa Maria de Palautordera (Catalonia, NE Spain). The authors of this study each made specific contributions to conceptualization, data collection, analysis, and/or writing and revising the manuscript. The primary data analyzed are described by Ward et al. (2019) and available in Ward (2019). Results of analyses completed in this study are available in Ward (2021). The authors declare no conflicts of interest. Any use of trade, firm, or product names is for descriptive purposes only and does not imply endorsement by the US government. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors.

## References

- Abbott, B. W., Gruau, G., Zarnetske, J. P., Moatar, F., Barbe, L., Thomas, Z., et al. (2018). Unexpected spatial stability of water chemistry in headwater stream networks. *Ecology Letters*, 21(2), 296–308. <https://doi.org/10.1111/ele.12897>
- Bernhardt, E. S., Blaszcak, J. R., Ficken, C. D., Fork, M. L., Kaiser, K. E., & Seybold, E. C. (2017). Control Points in Ecosystems: Moving Beyond the Hot Spot Hot Moment Concept. *Ecosystems*, 20(4), 665–682. <https://doi.org/10.1007/s10021-016-0103-y>
- Boulton, A. J., Harvey, M., & Proctor, H. (2004). Of spates and species: responses by interstitial water mites to simulated spates in a subtropical Australian river. *Exp Appl Acarol*, 34(1–2), 149–169.
- Briggs, MA, FD Day-Lewis, Zarnetske, JP, and JW Harvey (2015) A physical explanation for the development of redox microzones in hyporheic flow. *Geophysical Research Letters*, 42, doi: 10.1002/2015GL064200.
- Burt, T. P., & McDonnell, J. J. (2015). Whither field hydrology? The need for discovery science and outrageous hydrological hypotheses. *Water Resources Research*, 51. [https://doi.org/10.1016/0022-1694\(68\)90080-2](https://doi.org/10.1016/0022-1694(68)90080-2)
- Byrne, P., Wood, P. J., & Reid, I. (2012). The Impairment of River Systems by Metal Mine Contamination: A Review Including Remediation Options. *Critical Reviews in Environmental Science and Technology*, 42(19), 2017–2077. <https://doi.org/10.1080/10643389.2011.574103>
- Cardenas, M. B. (2008). Surface water-groundwater interface geomorphology leads to scaling of residence times. *Geophys. Res. Lett*, 35.
- Cotrufo, M. F., Wallenstein, M. D., Boot, C. M., Deneff, K., & Paul, E. (2013). The Microbial Efficiency-Matrix Stabilization (MEMS) framework integrates plant litter decomposition with soil organic matter stabilization: do labile plant inputs form stable soil organic matter? *Global Change Biology*, 19(4), 988–995. <https://doi.org/10.1111/GCB.12113>
- Czuba, J. A., David, S. R., Edmonds, D. A., & Ward, A. S. (2019). Dynamics of Surface-Water Connectivity in a Low-Gradient Meandering River Floodplain. *Water Resources Research*, 55(3). <https://doi.org/10.1029/2018WR023527>
- Danczak, R. E., Chu, R. K., Fansler, S. J., Goldman, A. E., Graham, E. B., Tfaily, M. M., et al. (2020). Using metacommunity ecology to understand environmental metabolomes. *Nature Communications* 2020 11:1, 11(1), 1–16. <https://doi.org/10.1038/s41467-020-19989-y>
- Deligne, N. I., McKay, D., Conrey, R. M., Grant, G. E., Johnson, E. R., O'Connor, J., & Sweeney, K. (2017). Field-trip guide to mafic volcanism of the Cascade Range in Central Oregon—A volcanic, tectonic, hydrologic, and geomorphic journey. *Scientific Investigations Report*, 110. <https://doi.org/10.3133/sir20175022H>
- Dyrness, C. T. (1969). Hydrologic properties of soils on three small watersheds in the western Cascades of Oregon. *USDA FOREST SERV RES NOTE PNW-111, SEP 1969. 17 P.*
- Fisher, S. G., Grimm, N. B., Martens, E., Holmes, R. M., & Jr., J. B. J. (1998). Material Spiraling in Stream Corridors: A Telescoping Ecosystem Model. *Ecosystems*, 1(1), 19–34. <https://doi.org/10.1007/s100219900003>
- Fondi, M., Karkman, A., Tamminen, M. V., Bosi, E., Virta, M., Fani, R., et al. (2016). “Every Gene Is Everywhere but the Environment Selects”: Global Geolocalization of Gene Sharing in Environmental Samples through Network Analysis. *Genome Biology and Evolution*, 8(5), 1388. <https://doi.org/10.1093/GBE/EVW077>

- Frissell, C. A., Liss, W. J., Warren, C. E., & Hurley, M. D. (1986). A hierarchical framework for stream habitat classification: Viewing streams in a watershed context. *Environmental Management*, 10(2), 199–214.
- Giraldo, L., Palacio, C., & Aguirre, N. (2014). Temporal Variation of the Extracellular Enzymatic Activity (EEA): Case of Study : Aburra-Medellín River, in the Valle de Aburra in Medellin, Antioquia, Colombia. *International Journal of Environmental Protection*, 4(5), 58–67.
- Graham, E. B., Crump, A. R., Kennedy, D. W., Arntzen, E., Fansler, S., Purvine, S. O., et al. (2018). Multi 'omics comparison reveals metabolome biochemistry, not microbiome composition or gene expression, corresponds to elevated biogeochemical function in the hyporheic zone. *Science of the Total Environment*, 642, 742–753. <https://doi.org/10.1016/j.scitotenv.2018.05.256>
- Gregory, K. J. (2006). The human role in changing river channels. *Geomorphology*, 79(3–4), 172–191. <https://doi.org/10.1016/j.geomorph.2006.06.018>
- Hampton, TB., JP Zarnetske, MA Briggs, F MahmoodPoor Dehkordy, K Singha, FD Day-Lewis, JW Harvey, S Roy Chowdhury and JW Lane. (2020) Experimental shifts of hydrologic residence time in a sandy urban stream sediment–water interface alter nitrate removal and nitrous oxide fluxes. *Biogeochemistry* 149, 195–219. <https://doi.org/10.1007/s10533-020-00674-7>.
- Hampton, TB, JP Zarnetske, MA Briggs, K Singha, JW Harvey, FD Day-Lewis, F MahmoodPoor Dehkordy, and JW Lane (2019) Residence time controls the fate of nitrogen in flow-through lakebed sediments. *JGR-Biogeosciences*, 124, 689– 707. <https://doi.org/10.1029/2018JG004741>
- Harman, C. J., Ward, A. S., & Ball, A. (2016). How does reach-scale stream-hyporheic transport vary with discharge? Insights fromrSAS analysis of sequential tracer injections in a headwater mountain stream. *Water Resources Research*, 52, 7130–7150. <https://doi.org/10.1002/2016WR018832>.Received
- Harvey, J. W., & Gooseff, M. N. (2015). River corridor science: Hydrologic exchange and ecological consequences from bedforms to basins. *Water Resources Research*, 51, 6893–6922. <https://doi.org/10.1002/2015WR017617>
- Heffernan, J. B., Soranno, P. A., Angilletta, M. J., Buckley, L. B., Gruner, D. S., Keitt, T. H., et al. (2014). Macrosystems ecology: understanding ecological patterns and processes at continental scales. *Frontiers in Ecology and the Environment*, 12(1), 5–14. <https://doi.org/10.1890/130017>
- Hill, B. H., McCormick, F. H., Harvey, B. C., Johnson, S. L., Warren, M. L., & Elonen, C. M. (2010). Microbial enzyme activity, nutrient uptake and nutrient limitation in forested streams. *Freshwater Biology*, 55(5), 1005–1019. <https://doi.org/10.1111/J.1365-2427.2009.02337.X>
- Isaak, D. J., Peterson, E. E., Ver Hoef, J. M., Wenger, S. J., Falke, J. A., Torgersen, C. E., et al. (2014). Applications of spatial statistical network models to stream data. *Wiley Interdisciplinary Reviews: Water*, 1(3), 277–294. <https://doi.org/10.1002/wat2.1023>
- Jefferson, A., Grant, G. E., & Lewis, S. L. (2004). A River Runs Underneath It: Geological Control of Spring and Channel Systems and Management Implications, Cascade Range, Oregon. In *Advancing the Fundamental Sciences Proceedings of the Forest Service: Proceedings of the Forest Service National Earth Sciences Conference* (Vol. 1, pp. 18–22).
- Jerolmack, D. J., & Paola, C. (2010). Shredding of environmental signals by sediment transport.



- Geophysical Research Letters*, 37(19), 1–5. <https://doi.org/10.1029/2010GL044638>
- Knapp, J. L. A., González-Pinzón, R., Drummond, J. D., Larsen, L. G., Cirpka, O. A., & Harvey, J. W. (2017). Tracer-based characterization of hyporheic exchange and benthic biolayers in streams. *Water Resources Res*, 53, 1575–1594. <https://doi.org/10.1002/2016WR019393>
- Krause, S., Hannah, D. M., Fleckenstein, J. H., Heppell, C. M., Kaeser, D. H., Pickup, R., et al. (2011). Inter-disciplinary perspectives on processes in the hyporheic zone. *Ecohydrology*, 4(4), 481–499.
- Krause, S., Lewandowski, J., Grimm, N. B., Hannah, D. M., Pinay, G., McDonald, K., et al. (2017). Ecohydrological interfaces as hot spots of ecosystem processes. *Water Resources Research*, 53(8), 6359–6376. <https://doi.org/10.1002/2016WR019516>
- Langbein, W. B., & Leopold, L. B. (1966). *River meanders - theory of minimum variance*.
- Lee-Cullin, J. A., Zarnetske, J. P., Ruhala, S. S., & Plont, S. (2018). Toward measuring biogeochemistry within the stream-groundwater interface at the network scale: An initial assessment of two spatial sampling strategies. *Limnology and Oceanography: Methods*, 16(11), 722–733. <https://doi.org/10.1002/lom3.10277>
- Leopold, L. B., Wolman, M. G., & Miller, J. P. (1964). *Fluvial Processes in Geomorphology*. Dover Publications.
- Leopold, L. B., & Langbein, W. B. (1962). *The Concept of Entropy in Landscape Evolution*.
- Li, L., Sullivan, P. L., Benettin, P., Cirpka, O. A., Bishop, K., Brantley, S. L., et al. (2021). Toward catchment hydro-biogeochemical theories. *Wiley Interdisciplinary Reviews: Water*, 8(1), e1495. <https://doi.org/10.1002/wat2.1495>
- Liébault, F., & Piégay, H. (2002). Causes of 20th century channel narrowing in mountain and piedmont rivers of southeastern France. *Earth Surface Processes and Landforms*, 27(4), 425–444. <https://doi.org/10.1002/esp.328>
- Lovett, G. M., Burns, D. A., Driscoll, C. T., Jenkins, J. C., Mitchell, M. J., Rustad, L., et al. (2007). Who needs environmental monitoring? *Frontiers in Ecology and the Environment*, 5(5), 253–260. [https://doi.org/10.1890/1540-9295\(2007\)5\[253:WNEM\]2.0.CO;2](https://doi.org/10.1890/1540-9295(2007)5[253:WNEM]2.0.CO;2)
- Martin, P. Y., & Turner, B. A. (1986). Grounded Theory and Organizational Research. *The Journal of Applied Behavioral Science*, 22(2), 141–157. <https://doi.org/10.1177/002188638602200207>
- McCluney, K. E., Poff, N. L., Palmer, M. A., Thorp, J. H., Poole, G. C., Williams, B. S., et al. (2014). Riverine macrosystems ecology: sensitivity, resistance, and resilience of whole river basins with human alterations. *Frontiers in Ecology and the Environment*, 12(1), 48–58. <https://doi.org/10.1890/120367>
- McGuire, K. J., Torgersen, C. E., Likens, G. E., Buso, D. C., Lowe, W. H., & Bailey, S. W. (2014). Network analysis reveals multiscale controls on streamwater chemistry. *Proceedings of the National Academy of Sciences of the United States of America*, 111(19), 7030–7035. <https://doi.org/10.1073/pnas.1404820111>
- Minshall, G. W., Petersen, R. C., Cummins, K. W., Bott, T. L., Sedell, J. R., Cushing, C. E., & Vannote, R. L. (1983). Interbiome Comparison of Stream Ecosystem Dynamics. *Ecological Monographs*, 53(1), 1–25. <https://doi.org/10.2307/1942585>
- Nikora, V. I., & Hicks, D. M. (1997). Scaling Relationships for Sand Wave Development in Unidirectional Flow. *Journal of Hydraulic Engineering*, 123(12), 1152–1156. [https://doi.org/10.1061/\(asce\)0733-9429\(1997\)123:12\(1152\)](https://doi.org/10.1061/(asce)0733-9429(1997)123:12(1152))
- O'Malley, M.A. (2008). “Everything is everywhere: but the environment selects”: ubiquitous distribution and ecological determinism in microbial biogeography. *Studies in History and*

*Philosophy of Biological and Biomedical Sciences*, 39(3), 314–325.

<https://doi.org/10.1016/J.SHPSC.2008.06.005>

- Payn, R. A., Gooseff, M. N., McGlynn, B. L., Bencala, K. E., & Wondzell, S. M. (2009). Channel water balance and exchange with subsurface flow along a mountain headwater stream in Montana, United States. *Water Resources Research*, 45.
- Pinay, G., Peiffer, S., De Dreuzy, J. R., Krause, S., Hannah, D. M., Fleckenstein, J. H., et al. (2015). Upscaling Nitrogen Removal Capacity from Local Hotspots to Low Stream Orders' Drainage Basins. *Ecosystems*, 18(6), 1101–1120. <https://doi.org/10.1007/s10021-015-9878-5>
- Pringle, C. M., Naiman, R. J., Bretschko, G., Karr, J. R., Oswood, M. W., Webster, J. R., et al. (1988). Patch Dynamics in Lotic Systems: The Stream as a Mosaic. *Journal of the North American Benthological Society*, 7(4), 503–524. <https://doi.org/10.2307/1467303>
- Rana, S. M. M., Scott, D. T., & Hester, E. T. (2017). Effects of in-stream structures and channel flow rate variation on transient storage. *Journal of Hydrology*, 548, 157–169. <https://doi.org/10.1016/j.jhydrol.2017.02.049>
- Robertson, A. D., Paustian, K., Ogle, S., Wallenstein, M. D., Lugato, E., & Francesca Cotrufo, M. (2019). Unifying soil organic matter formation and persistence frameworks: The MEMS model. *Biogeosciences*, 16(6), 1225–1248. <https://doi.org/10.5194/BG-16-1225-2019>
- Rodríguez-Iturbe, I., & Rinaldo, A. (1997). *Fractal River Basins: Chance and Self-Organization*. Cambridge, UK: Cambridge University Press.
- Santschi, P. H., Presley, B. J., Wade, T. L., Garcia-Romero, B., & Baskaran, M. (2001). Historical contamination of PAHs, PCBs, DDTs, and heavy metals in Mississippi River Delta, Galveston Bay and Tampa Bay sediment cores. *Marine Environmental Research*, 52(1), 51–79. [https://doi.org/10.1016/S0141-1136\(00\)00260-9](https://doi.org/10.1016/S0141-1136(00)00260-9)
- Sinsabaugh, R. L., Findlay, S., Franchini, P., & Fischer, D. (1997). Enzymatic analysis of riverine bacterioplankton production. *Limnology and Oceanography*, 42(1), 29–38. <https://doi.org/10.4319/LO.1997.42.1.0029>
- Sinsabaugh, R. L., Findlay, S., Franchini, P., & Fischer, D. (1997). Enzymatic analysis of riverine bacterioplankton production. *Limnology and Oceanography*, 42(1), 29–38. <https://doi.org/10.4319/LO.1997.42.1.0029>
- Sinsabaugh, R. L., & Shah, J. J. F. (2012). Ecoenzymatic Stoichiometry and Ecological Theory. <http://Dx.Doi.Org/10.1146/Annurev-Ecolsys-071112-124414>, 43, 313–343. <https://doi.org/10.1146/ANNUREV-ECOLSYS-071112-124414>
- Smidt, S. J., Cullin, J. A., Ward, A. S., Robinson, J., Zimmer, M. A., Lutz, L. K., & Endreny, T. A. (2015). A Comparison of Hyporheic Transport at a Cross-Vane Structure and Natural Riffle. *Ground Water*, 53(6), 859–871. <https://doi.org/10.1111/gwat.12288>
- Sollins, P., Cromack, K., Corison, F. M. M., Waring, R. H., & Harr, R. D. (1981). Changes in Nitrogen Cycling at an Old-Growth Douglas-fir Site After Disturbance. *Journal of Environmental Quality*, 10(1), 37–42. <https://doi.org/10.2134/JEQ1981.00472425001000010007X>
- Stegen, J. C., & Goldman, A. E. (2018). WHONDRS: a Community Resource for Studying Dynamic River Corridors. *MSystems*, 3(5), 151–169. <https://doi.org/10.1128/msystems.00151-18>
- Stegen, J. C., Bottos, E. M., & Jansson, J. K. (2018). A unified conceptual framework for prediction and control of microbiomes. *Current Opinion in Microbiology*, 44(July), 20–27. <https://doi.org/10.1016/j.mib.2018.06.002>

- Stegen, J. C., Johnson, T., Fredrickson, J. K., Wilkins, M. J., Konopka, A. E., Nelson, W. C., et al. (2018). Influences of organic carbon speciation on hyporheic corridor biogeochemistry and microbial ecology. *Nature Communications*, 9(1), 1–11. <https://doi.org/10.1038/s41467-018-03572-7>
- Strauss, A., & Corbin, J. (1994). Grounded theory methodology: An overview. In N. Denzin & Y. Lincoln (Eds.), *Handbook of qualitative research* (pp. 273–285). Sage Publications, Inc.
- Swanson, F. J., & James, M. E. (1975). *Geology and geomorphology of the H.J. Andrews Experimental Forest, western Cascades, Oregon*. Portland, OR.
- Swanson, F. J., & Jones, J. A. (2002). Geomorphology and hydrology of the H.J. Andrews Experimental Forest, Blue River, Oregon. In *Field guide to geologic processes in Cascadia*.
- Tank, J. L., Rosi-Marshall, E. J., Baker, M. A., & Hall, R. O. (2008). Are rivers just big streams? A pulse method to quantify nitrogen demand in a large river. *Ecology*, 89(10), 2935–2945.
- Tiegs, S. D., Costello, D. M., Isken, M. W., Woodward, G., McIntyre, P. B., Gessner, M. O., et al. (2019). Global patterns and drivers of ecosystem functioning in rivers and riparian zones. *Science Advances*, 5(1), eaav0486. <https://doi.org/10.1126/SCIADV.AAV0486>
- Triska, F. J., Sedell, J. R., Cromack, K., Gregory, S. V., & McCorison, F. M. (1984). Nitrogen Budget for a Small Coniferous Forest Stream. *Ecological Monographs*, 54(1), 119–140. <https://doi.org/10.2307/1942458>
- Turnbull, L., Hütt, M. T., Ioannides, A. A., Kininmonth, S., Poepl, R., Tockner, K., et al. (2018, December 1). Connectivity and complex systems: learning from a multi-disciplinary perspective. *Applied Network Science*. Springer. <https://doi.org/10.1007/s41109-018-0067-2>
- Valett, H. M., Morrice, J. A., Dahm, C. N., & Campana, M. E. (1996). Parent lithology, surface-groundwater exchange, and nitrate retention in headwater streams. *Limnology and Oceanography*, 333–345.
- Vannote, R. L., Minshall, G. W., Cummins, K. W., Sedell, J. R., & Cushing, C. E. (1980). The River Continuum Concept. *Canadian Journal of Fisheries and Aquatic Sciences*, 37, 130–137.
- Ver Hoef, J. M., Peterson, E., & Theobald, D. (2006). Spatial statistical models that use flow and stream distance. *Environmental and Ecological Statistics*, 13(4), 449–464. <https://doi.org/10.1007/s10651-006-0022-8>
- Walling, D. E., & Fang, D. (2003). Recent trends in the suspended sediment loads of the world's rivers. *Global and Planetary Change*, 39(1–2), 111–126. [https://doi.org/10.1016/S0921-8181\(03\)00020-1](https://doi.org/10.1016/S0921-8181(03)00020-1)
- Wallis, I., Prommer, H., Berg, M., Siade, A. J., Sun, J., & Kipfer, R. (2020). The river–groundwater interface as a hotspot for arsenic release. *Nature Geoscience*, 13(4), 288–295. <https://doi.org/10.1038/s41561-020-0557-6>
- Ward, A. S. (2015). The evolution and state of interdisciplinary hyporheic research. *Wiley Interdisciplinary Reviews: Water*, 3(1), 83–103. <https://doi.org/10.1002/wat2.1120>
- Ward, A. S. (2019). ESSD, 2019 - Data Collection. <https://doi.org/10.5194/essd-11-1-2019>
- Ward, A. S. (2021). Supporting data for Ward et al., (In Review) Advancing river corridor science beyond disciplinary boundaries with an inductive approach to hypothesis generation, HydroShare, Accessed 6-May-2021. <http://www.hydroshare.org/resource/de6d92d314354ea6819157818669fc59>
- Ward, A. S., & Packman, A. I. (2019). Advancing our predictive understanding of river corridor exchange. *Wiley Interdisciplinary Reviews: Water*, 6(1), e1327. <https://doi.org/10.1002/wat2.1327>

- Ward, A. S., Zarnetske, J. P., Baranov, V., Blaen, P. J., Brekenfeld, N., Chu, R., et al. (2019). Co-located contemporaneous mapping of morphological, hydrological, chemical, and biological conditions in a 5th-order mountain stream network, Oregon, USA. *Earth System Science Data*, 11(4). <https://doi.org/10.5194/essd-11-1567-2019>
- Ward, A. S., Wondzell, S. M., Schmadel, N. M., Herzog, S., Zarnetske, J. P., Baranov, V., et al. (2019). Spatial and temporal variation in river corridor exchange across a 5th order mountain stream network. *Hydrology and Earth System Sciences Discussions*, (April), 1–39. <https://doi.org/10.5194/hess-2019-108>
- Waring, B. G., Sulman, B. N., Reed, S., Smith, A. P., Averill, C., Creamer, C. A., et al. (2020). From pools to flow: The PROMISE framework for new insights on soil carbon cycling in a changing world. *Global Change Biology*, 26(12), 6631–6643. <https://doi.org/10.1111/GCB.15365>
- Williams, C. J., Scott, A. B., Wilson, H. F., & Xenopoulos, M. A. (2011). Effects of land use on water column bacterial activity and enzyme stoichiometry in stream ecosystems. *Aquatic Sciences* 2011 74:3, 74(3), 483–494. <https://doi.org/10.1007/S00027-011-0242-3>
- Williams, C. J., Yamashita, Y., Wilson, H. F., Jaffé, R., & Xenopoulos, M. A. (2010). Unraveling the role of land use and microbial activity in shaping dissolved organic matter characteristics in stream ecosystems. *Limnology and Oceanography*, 55(3), 1159–1171. <https://doi.org/10.4319/LO.2010.55.3.1159>
- Wisnoski, N. I., & Lennon, J. T. (2021). Microbial community assembly in a multi-layer dendritic metacommunity. *Oecologia*, 195(1), 13–24. <https://doi.org/10.1007/s00442-020-04767-w>
- Wisnoski, N. I., Muscarella, M. E., Larsen, M. L., Peralta, A. L., & Lennon, J. T. (2020). Metabolic insight into bacterial community assembly across ecosystem boundaries. *Ecology*, 101(4), e02968. <https://doi.org/10.1002/ECY.2968>
- Wit, R. De, & Bouvier, T. (2006). ‘Everything is everywhere, but, the environment selects’; what did Baas Becking and Beijerinck really say? *Environmental Microbiology*, 8(4), 755–758. <https://doi.org/10.1111/J.1462-2920.2006.01017.X>
- Wohl, E. (2005). Compromised Rivers: Understanding Historical Human Impacts on Rivers in the Context of Restoration. *Ecology and Society*, 10(2), 2.
- Wondzell, S. M., & Gooseff, M. N. (2014). Geomorphic Controls on Hyporheic Exchange Across Scales: Watersheds to Particles. In J. Schroder & E. Wohl (Eds.), *Treatise on Geomorphology* (Vol. 9, pp. 203–218). San Diego, CA: Academic Press.
- Wood, P. J., Boulton, A. J., Little, S., & Stubbington, R. (2010). Is the hyporheic zone a refugium for aquatic macroinvertebrates during severe low flow conditions? *Fundamental and Applied Limnology / Archiv Für Hydrobiologie*, 176(4), 377–390. <https://doi.org/10.1127/1863-9135/2010/0176-0377>
- Wörman, A., Packman, A. I., Marklund, L., Harvey, J. W., & Stone, S. H. (2007). Fractal topography and subsurface water flows from fluvial bedforms to the continental shield. *Geophysical Research Letters*, 34(7), 1–5. <https://doi.org/10.1029/2007GL029426>
- Wu, L., Singh, T., Gomez-Velez, J., Nützmann, G., Wörman, A., Krause, S., & Lewandowski, J. (2018). Impact of Dynamically Changing Discharge on Hyporheic Exchange Processes Under Gaining and Losing Groundwater Conditions. *Water Resources Research*, 54(12), 10,076–10,093. <https://doi.org/10.1029/2018WR023185>
- Yoder, L., Ward, A. S., Spak, S., & Dalrymple, K. (2020). Local Government Perspectives on Collaborative Governance: A Comparative Analysis of Iowa’s Watershed Management

Authorities. *Policy Studies Journal*. <https://doi.org/10.1111/psj.12389>

Zhou, C., Liu, Y., Liu, C., Liu, Y., & Tfaily, M. M. (2019). Compositional changes of dissolved organic carbon during its dynamic desorption from hyporheic zone sediments. *Science of the Total Environment*, 658, 16–23. <https://doi.org/10.1016/j.scitotenv.2018.12.189>