



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

Journal:	<i>Hydrological Processes</i>
Manuscript ID	HYP-21-0678.R1
Wiley - Manuscript type:	Special Issue Paper (Direct Via EEO)
Date Submitted by the Author:	n/a
Complete List of Authors:	<p>Ward, Adam; Indiana University, School of Public and Environmental Affairs Packman, Aaron; Northwestern University , CEE; Bernal, Susana; Centre of Advanced Studies in Blanes, Geodynamics and Biodiversity Group Brekenfeld, Nicolai; University of Birmingham Edgbaston Campus, Birmingham Institute of Forest Research Drummond, Jennifer; University of Birmingham Edgbaston Campus, School of Geography, Earth and Environmental Sciences Graham, Emily; Pacific Northwest National Laboratory, EMSL Hannah, David; University of Birmingham, Geography, Earth and Environmental Sciences Klaar, Megan; University of Leeds, School of Geography Krause, Stefan; University of Birmingham, School of Geography, Earth and Environmental Sciences Kurz, Marie; Drexel University, Academy of Natural Sciences Li, Angang; Northwestern University, Civil and Environmental Engineering Lupon, Anna; Centre of Advanced Studies in Blanes, Geodynamics and Biodiversity Group Mao, Feng; Cardiff University, School of Earth and Environmental Sciences Martí, Eugenia; Centre d'Estudis Avançats de Blanes (CSIC), Aquatic Biogeochemistry Ouellet, Valerie; NOAA, Fisheries Royer, Todd; Indiana University System, SPEA Stegen, James; Pacific Northwest National Laboratory, Earth and Biological Sciences Directorate Zarnetske, Jay; MICHIGAN STATE UNIVERSITY, Dept. of Earth and Environmental Science</p>
Keywords:	river corridor, stream corridor, machine learning, inductive, scientific method

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



1
2
3 **Advancing river corridor science beyond disciplinary boundaries with an inductive**
4 **approach to catalyze hypothesis generation**
5

6 *Submitted for publication in Hydrological Processes – Special issue – Data Science Applications*
7 *in Hydrology*
8
9

10 **Authors:**

11 Adam S. Ward¹, Aaron Packman², Susana Bernal³, Nicolai Brekenfeld⁴, Jen Drummond⁴, Emily
12 Graham⁵, David M. Hannah⁴, Megan Klaar⁶, Stefan Krause⁴, Marie Kurz⁷, Angang Li², Anna
13 Lupon³, Feng Mao⁸, M. Eugènia Martí Roca³, Valerie Ouellet⁴, Todd V. Royer¹, James C.
14 Stegen⁵, Jay P. Zarnetske⁹
15

16
17 ¹ O'Neill School of Public and Environmental Affairs, Indiana University, Bloomington,
18 Indiana, USA

19 ² Department of Civil and Environmental Engineering, Northwestern University, Evanston,
20 Illinois, USA

21 ³ Integrative Freshwater Ecology Group, Centre for Advanced Studies of Blanes (CEAB-
22 CSIC), Blanes, Spain

23 ⁴ School of Geography, Earth & Environmental Sciences, University of Birmingham,
24 Edgbaston, Birmingham, B15 2TT, UK

25 ⁵ Earth and Biological Sciences Directorate, Pacific Northwest National Laboratory, Richland,
26 Washington, USA

27 ⁶ School of Geography, School of Earth and Environment, University of Leeds, Woodhouse,
28 Leeds LS2 9JT, United Kingdom

29 ⁷ The Academy of Natural Sciences of Drexel University, Philadelphia, Pennsylvania, USA

30 ⁸ School of Earth and Environmental Sciences, Cardiff University, Building, Park Place,
31 Cardiff, CF10 3AT, United Kingdom

32 ⁹ Department of Earth and Environmental Sciences, Michigan State University, East Lansing,
33 Michigan, USA
34
35
36

37
38 **Corresponding author:**

39 Adam S. Ward
40 O'Neill School of Public and Environmental Affairs
41 Indiana University
42 418 MSB-II
43 Bloomington, IN 47405
44

45
46 Email: adamward@indiana.edu
47 Phone: 812-865-4820
48

49 **Running head:** Inductive hypothesis generation using data science
50

51 **Key words:** river corridor, stream corridor, machine learning, inductive, scientific method
52
53
54
55
56
57
58

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 5th 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

1
2
3 approach to science, which bases measurements on a highly targeted set of causal mechanisms to
4 be tested at a limited range of locations and scales. With the emergence of new experimental and
5 data science techniques, the time has come to expand existing conceptual models for river
6 corridors via approaches that generate more integrative knowledge commensurate with the
7 reality of of river corridors as complex dynamic systems. We posit that unified understanding
8 must be derived from a combination of *deductive* science and *inductive* approaches that identify
9 process interactions and couplings that emerge from the data themselves. We suggest that river
10 corridor science can benefit from inductive approaches that generate hypotheses and eventually
11 theories from empirical studies, an approach successfully applied in other disciplines (Martin &
12 Turner, 1986; Strauss & Corbin, 1994; e.g., Turnbull et al., 2018).

13
14
15
16
17
18
19
20
21
22 A unifying framework is required to organize and synthesize our understanding of river corridors
23 and advance scientific understanding of the drivers and controls of their functioning. Stegen et al.
24 (2018) propose one such model for microbial ecology, where the resultant ecosystem functions
25 and services are explained by the relationships linking internal dynamics, external forcing, and
26 historical contingencies. The principles of Stegen et al.'s conceptual framework are similar to
27 other existing conceptualizations of river corridors that have been developed by other disciplines.
28 First, external forcing describes the role of factors extrinsic to the river corridor that shape its
29 structure and function. For river corridors, this primarily means the larger spatial scale and
30 longer temporal scale elements that are functionally decoupled (e.g., static or slowly-varying)
31 relative to a process of interest. Studies with data collection spanning gradients in land use,
32 geologic setting, climate, network position, or other factors that are considered to be extrinsic
33 typically use geospatial and statistical approaches to describe patterns and trends (e.g., McGuire
34 et al., 2014), while variation around spatially structured trends is often interpreted as random
35 noise from structural heterogeneity and/or unstudied, smaller-scale processes (Abbott et al.,
36 2018). Next, internal dynamics are the interacting processes within the river corridor that give
37 rise to observed functions of interest at a given location. Conceptual models based on this
38 approach to river corridor science include hot spots and hot moments (Krause et al., 2011, 2017;
39 Wallis et al., 2020), control points (Bernhardt et al., 2017), and patch dynamics (Pringle et al.,
40 1988). River corridor dynamics are commonly studied through detailed observations at a
41 relatively limited spatial scale, which is restricted in an attempt to characterize local feedbacks
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 between mechanisms. These approaches often lack sufficient spatial resolution to enable
4 confident application of geostatistical approaches, and may not reliably support assessments of
5 system dynamics (e.g., Lee-Cullin et al., 2018). Longer-term dynamics are often considered as
6 historical contingencies: the biotic and abiotic histories or antecedent conditions that lead to the
7 present characteristics of the river corridor and affect its response to future perturbations.
8
9
10 Examples of river corridor studies that incorporate historical contingencies include perturbation-
11 response dynamics, commonly associated with floods (Czuba et al., 2019; Wu et al., 2018),
12 droughts (Boulton et al., 2004; Wood et al., 2010), or restoration activities (Rana et al., 2017;
13 Smidt et al., 2015), and large-scale historical perturbations such as land development (Liébault &
14 Piégay, 2002; Walling & Fang, 2003; Wohl, 2005), river regulation (Gregory, 2006), and
15 contamination (Byrne et al., 2012; Santschi et al., 2001). Such studies often involve little to no
16 replication and may be biased towards response variables that change rapidly relative to
17 processes that are quasi-steady over the timeframe of a given experiment.
18
19
20
21
22
23
24
25
26

27 While external forcing, internal dynamics, and historical contingencies have each been studied in
28 their own right, recent studies are beginning to integrate these concepts into holistic
29 understanding of river corridors. For example, Wisnoski and Lennon (2021) explicitly linked
30 localized heterogeneity to systematic spatial patterns along the network, revealing that the local
31 microbial assemblage in headwaters streams was controlled by local physical and chemical
32 conditions, but these local controls gave way to systemic organization from headwaters to larger
33 downstream rivers as the spatial scale of study increased. Such explicit consideration of local and
34 network scales is rare and still does not address historical contingencies. However, if done more
35 often and expanded to consider historical contingencies as a context for each replicate, this type
36 of systematic approach would allow assessment of the transition in dominant controls from local
37 heterogeneity (a reflection of internal dynamics) to larger-scale spatial organization (a reflection
38 of external drivers), the specific mechanisms of this transition, and the scale at which the
39 transition occurs, and how future multi-scale dynamics may depend on antecedent conditions (a
40 reflection of historical contingencies). Studies that have explicitly considered local
41 spatiotemporal dynamics as part of long-term system-wide functions have found strong
42 relationships between large-scale system structure, internal dynamics, and long-term emergent
43 outcomes in flow, sediment transport, and biogeochemistry (e.g., Fisher et al., 1998; Harvey &
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Gooseff, 2015; Krause et al., 2017; Pinay et al., 2015). The success of these studies demonstrates
4 our ability to identify a core set of transferable and scalable processes that govern river system
5 dynamics and unify seemingly disparate observations into holistic understanding of river
6 corridor services and functions.
7
8
9

10
11 Here we use objective data-oriented approaches to confront the challenge of how a discipline
12 organized around the classic deductive model of site-, scale-, and mechanism-specific
13 observations can systematically link the resulting fragmented information into system-level
14 understanding. Our aim is to identify couplings that span scales and disciplinary expertise in
15 absence of pre-existing conceptual models that would traditionally serve as the source of
16 hypotheses for deductive testing. We propose an inductive approach to data synthesis, serving as
17 a basis for the unconstrained generation of new and potentially unexpected relationships, each of
18 which may be explained by hypotheses that could subsequently be tested. To this end, we
19 analyze a novel large data set for a 5th order river basin (Ward, Zarnetske, et al., 2019) using
20 inductive approaches to generate a network of relationships that span traditional disciplinary
21 boundaries. The data set contains 157 variables with nearly 25,000 possible pairwise
22 relationships, making it infeasible to explore each potential relationship through the lens of
23 deductive inquiry. Further, the large degree of covariation in environmental conditions may
24 obscure underlying causal mechanisms, making it difficult to determine unique process
25 relationships and their controls. Thus, we pilot a machine learning approach that sieves and
26 categorizes information to identify non-obvious relationships that merit subsequent investigation.
27 We envision the apparent relationships generated by our approach as a suite of observations
28 around which hypotheses could be generated and subsequently tested with more traditional
29 approaches. In this way, we complement traditional approaches by highlighting observations that
30 may warrant hypotheses to be spun that explain causal pathways that novel, interdisciplinary,
31 and trans-scale to explain the apparent relationships. This allows us to synthesize complex,
32 multi-scale observations independent of any pre-conceived conceptual models and uncover novel
33 and exciting information about the structure and function of river corridors. We critically
34 evaluate the resultant relationships relative to existing knowledge, and provide two examples of
35 how these novel insights may motivate future research questions that inform a synthesis
36 approach to understanding of river corridors.
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

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 5th 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

1
2
3 subgroups representing distinct study domains that jointly characterize the structure, function,
4 and dynamics of the river corridor and consistent with the design of the field campaign. These
5 subgroups were: geologic setting (GEO), physical chemistry (PCHEM), bulk DOM
6 characterization (DOM), dissolved nutrients (NUTS), solute tracers (TRACER), metabolomics
7 (ICR), and macroinvertebrates (MACRO). A complete list of variables, subgroups, and summary
8 findings for each variable is presented in Table S1). The reduced data set totaled 157 unique
9 variables across the seven disciplinary subgroups and is the basis for all subsequent analysis in
10 this study.
11
12
13
14
15
16
17
18

19 ***2.2 Principal components analysis***

20 To identify major axes of (co)variation among measured variables, we performed a series of
21 principal component analyses (PCAs) using the rotated PCA approach. Independent PCAs were
22 performed first on the entire data set (all 157 variables) and subsequently on variables within
23 each subgroup. For each PCA, we focused on results from the first two components (PC1 and
24 PC2). We identified the most influential variables from each principal component as those with
25 loadings greater than 0.6 or less than -0.6 (hereafter ‘influential variables’) and interpreted the
26 variables aligned with each PC to describe the major axes of variation when possible.
27
28
29
30
31
32
33

34 ***2.3 Spatial structure of individual variables***

35 For each variable, we tested for spatial structure throughout the network by assessing the change
36 in variance as a function of distance between flow connected points, (i.e., a semivariogram; Ver
37 Hoef et al., 2006; Isaak et al., 2014; McGuire et al., 2014). This analysis identifies variables for
38 which variance is spatially uniform (i.e., no change in variance as a function of distance),
39 increases linearly (i.e., variance grows with distance), or that plateaus at some distance (a scale
40 cutoff). A uniform relationship indicates no structure (hereafter, unstructured variable), while
41 both linear and plateau relationships demonstrate spatial structure (hereafter, structured variable).
42 The linear models were only considered significant if the estimate of the slope was significantly
43 different from zero based on the 95% confidence interval for a linear model fit. The squared
44 differences were normalized (squared difference subtracted from the mean, followed by division
45 of the difference by the standard deviation) and binned (bin size of 30) before being fitted. An
46 exponential semivariogram function was considered for cases that exhibited scale cutoffs:
47
48
49
50
51
52
53
54
55
56
57
58
59
60

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

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

1
2
3 The PCA on all variables identified major axes of co-variation without regard to disciplinary
4 grouping. PC1 explained 20% of the total variance (Table 2A), and contained mainly variables
5 from the metabolomics subgroup, generally representing a gradient moving from terrestrially-
6 derived aromatic compounds that are more thermodynamically favorable for microbial
7 respiration to more microbially-derived compounds that are less thermodynamically favorable.
8 PC2 explained 17% of the total variance and contained variables from the geologic setting
9 subgroup, such as valley width and stream slope, showing marked gradients from headwaters to
10 downstream reaches. Taken together PC1 and PC2 suggest that sampling sites within the river
11 network are organized by organic matter chemistry and geology, which are themselves linked by
12 terrestrial vegetation and soils.
13
14
15
16
17
18
19
20
21

22 ***3.1.2 Principal component analysis on disciplinary subgroups***

23
24 PCAs were conducted on each subgroup to identify major axes of variation within individual
25 disciplinary perspectives. The first two PCs within each subgroup explain an average of 52% of
26 the within-group variance (median 46%, range 33-76%; Fig. 1A; Table 1). For physical
27 chemistry, we interpret PC1 as representing weathering rate (from high to low) and PC2 as
28 representing age of water (from high to low). For the geophysical setting, we interpret PC1 as
29 representing network position (from headwaters to larger rivers) and PC2 as representing
30 surficial geology. For nutrients, we interpret PC1 as representing enzymatic activity (low to
31 high) which is itself the inverse of dissolved inorganic nutrient availability, and PC2 represents
32 the accumulated organic matter in the shallow streambed. For metabolomics, we interpret PC1 as
33 reflecting gradients from terrestrially-derived aromatic compounds that are more
34 thermodynamically favorable for microbial respiration to more microbially-derived compounds
35 that are less thermodynamically favorable. The metabolomics PC2 is interpreted as a gradient
36 being dominated by products from organic matter degradation at one end and less-processed
37 terrestrially-derived organic matter at the other end. For bulk DOM, we interpret PC1 as
38 representing DOM quality from less to more humic or terrestrial in origin, and PC2 as
39 representing microbial and proteinaceous DOM (from more to less). For macroinvertebrates, we
40 interpret PC1 as representing richness (high to low) and PC2 as representing abundance (high to
41 low). For stream solute tracers, we interpret PC1 as representing short-term storage of tracers
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

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

For Peer Review

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

PCA on all variables						
	PC1			PC2		
	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
PCA on subgroups						
	PC1			PC2		
	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 5th 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

1
2
3 unconstrained in the longitudinal dimension until they reach the ocean. Still others have found
4 spatial structure in some parameters (e.g., in-stream solute concentrations) at scales that were
5 encapsulated within our study (e.g., McGuire et al., 2014), suggesting that finding of spatial
6 correlation lengths in one system or for one variable may not be universally transferable.
7
8
9

10
11
12 The fraction of influential variables with spatial structure was varied between subgroups (Fig.
13 1B, 1C), with 6 of 14 subgroup of PCs containing both structured and unstructured variables.
14 The largest proportion of spatially structured variables were in the TRACER subgroup (69%;
15 Fig. 1C), and the least were in the PCHEM subgroup (9.5%; Fig. 1C). The variables that appear
16 in the disciplinary subgroup PCs did not separate into distinct groups of structured vs.
17 unstructured variables. Instead, we found 44% of all influential variables were spatially
18 structured (23% in PC1 and 21% in PC2) compared to overall 36% of all variables exhibiting
19 spatial structure. All subgroups contained some structured influential variables except for
20 MACRO (Fig. 1B), where only unstructured variables were selected.
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

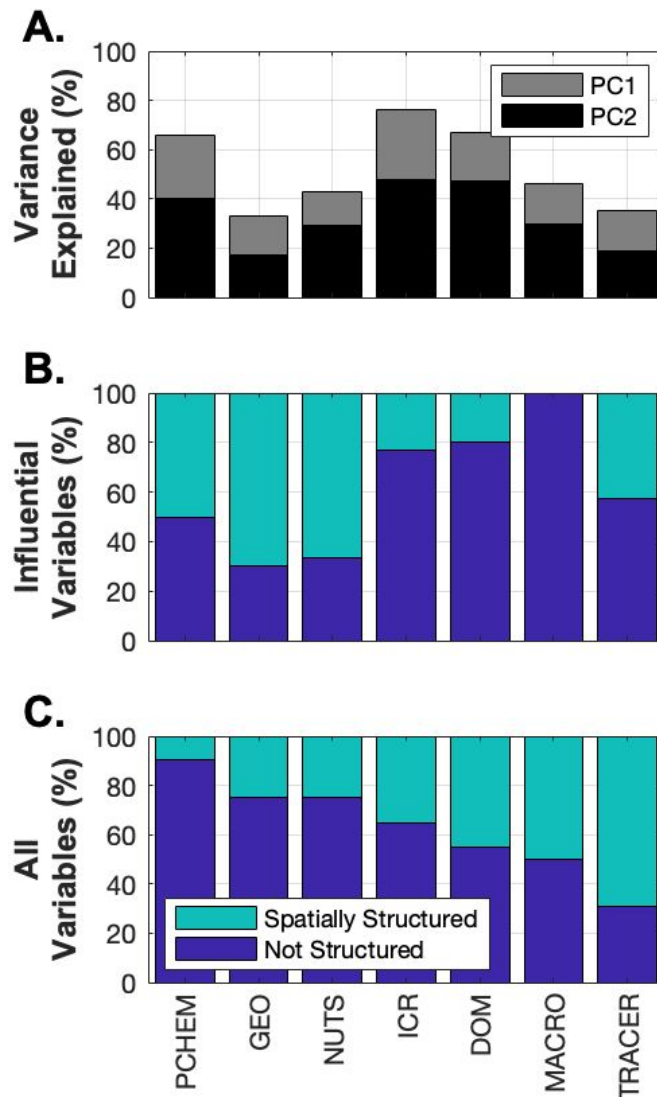


Fig. 1. (A) Variance in the Andrews river corridor data set explained by PC1 and PC2 for each expert subgroup. (B) Percentage of influential variables (i.e., the variables included in the first two PCs) that do and do not have spatial structure. (C) Percentage of all variables within each subgroup that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

1
2
3 in the underlying data, with an overall mean r^2 of 0.83 (median 0.94, range 0.00 - 1.00). SVMRs
4 for individual variables selected an average of 4.4 variables as predictors (median 4, range 1 to
5 10; Fig. S1), indicating that the relationships (i.e., statistical models) identified by the SVMRs
6 were reasonably parsimonious. Additionally, performance of the SVMRs built from the full
7 feature set was significantly better than those built from a random selection of features ($p_{ANOVA} =$
8 $1E-19$; $p_{KW} = 4E-29$), indicating SVMRs are selecting meaningful features and the associated
9 relationships are appropriate for further analysis. The models built for spatially structured
10 variables had an overall mean r^2 of 0.91 (median 0.97, range 0.08 - 1.00) compared to a mean r^2
11 of 0.78 for unstructured variables (median 0.90, range 0.00 - 1.00). Goodness of fit was also
12 statistically better for the spatially structured variables ($p = 0.008$; one-way ANOVA), indicating
13 that spatially structured variables were more accurately predicted (i.e., higher r^2) compared to
14 unstructured variables.
15
16
17
18
19
20
21
22
23
24

25 Of the 157 variables predicted, 22% (34 variables) are informed by only out-of-group variables
26 (i.e., variables from a different subgroup), and 11% (17 variables) are informed by only within-
27 group variables (i.e., variables in the same subgroup). Thus, 67% of variables (106 out of 157)
28 required both in-group and out-of-group information for optimal prediction by the SVMRs.
29 Moreover, out-of-group information dominates predictor selection, representing an average of
30 59% of variables selected (median 66%, range 0-100%; Fig. 2, Table S1). Spatially structured
31 variables represent an average of 27.3% of variables selected for individual SVMRs (Fig. S2A,
32 S2C). Across the 157 SVMRs constructed, 30% (47 variables) did not select any spatially
33 structured features. We found 3% of models (5 variables) selected only spatially structured
34 features, and the remaining 67% (105 variables) selected a combination of structured and
35 unstructured variables.
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

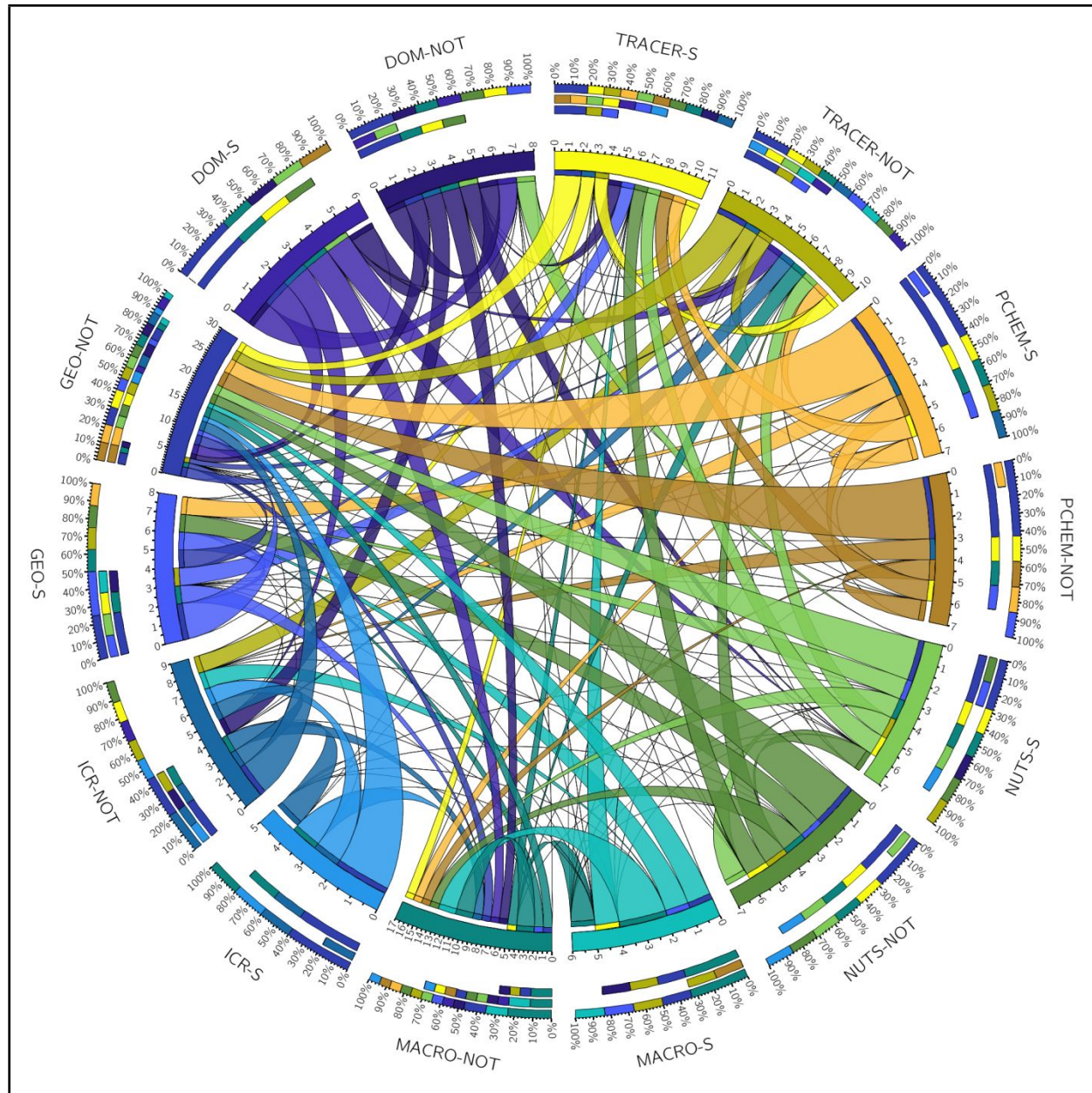


Fig. 2. Information flow within and among subgroups of variables commonly used as measures of river corridor dynamics based on the suite of SVMRs constructed for each variable (Section 3.3.1). The variables included in the 7 subgroups (PCHEM = physical chemistry; GEO = geologic setting; NUTS = nutrients; 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 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.

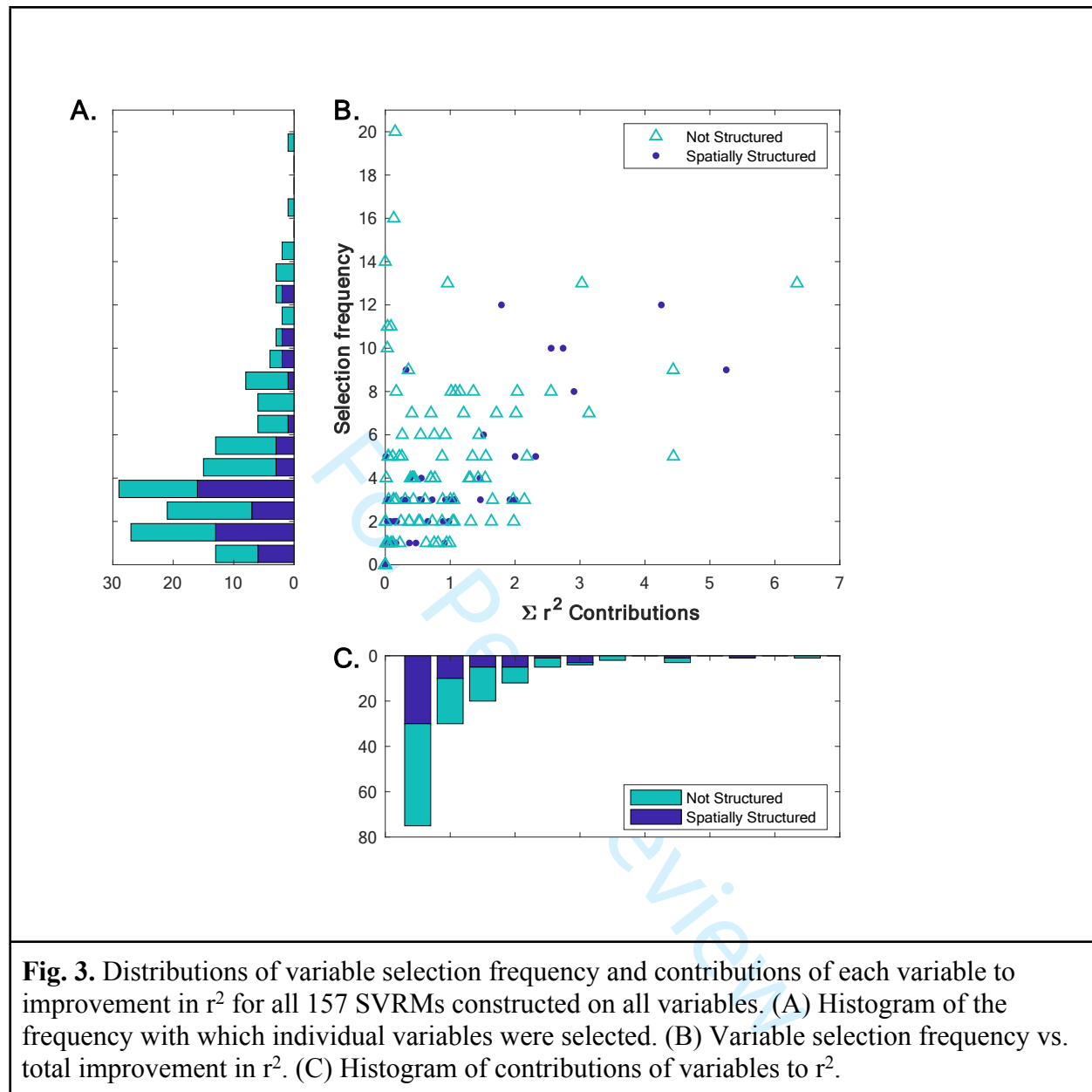


Fig. 3. Distributions of variable selection frequency and contributions of each variable to improvement in r^2 for all 157 SVRMs constructed on all variables. (A) Histogram of the frequency with which individual variables were selected. (B) Variable selection frequency vs. total improvement in r^2 . (C) Histogram of contributions of variables to r^2 .

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-

1
2
3 group variance is unsurprising given that PC1 and PC2 were successful in explaining an average
4 of 52% of variance within their group. However, we also found about 26% of variable
5 predictions (41 of 157) used only out-of-group PCs, and 118 variable predictions selected at least
6 one out-of-group PCs. Further, variables in each subgroup drew information from nearly every
7 other subgroup (see Table S1), These findings indicate that studies that are limited to one
8 discipline are unlikely to explain as much of the observed variance in the measured variables as
9 studies that intentionally span disciplinary boundaries, and that it is important for disciplinary
10 understanding to at least characterize the major attributes from other subgroups.
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

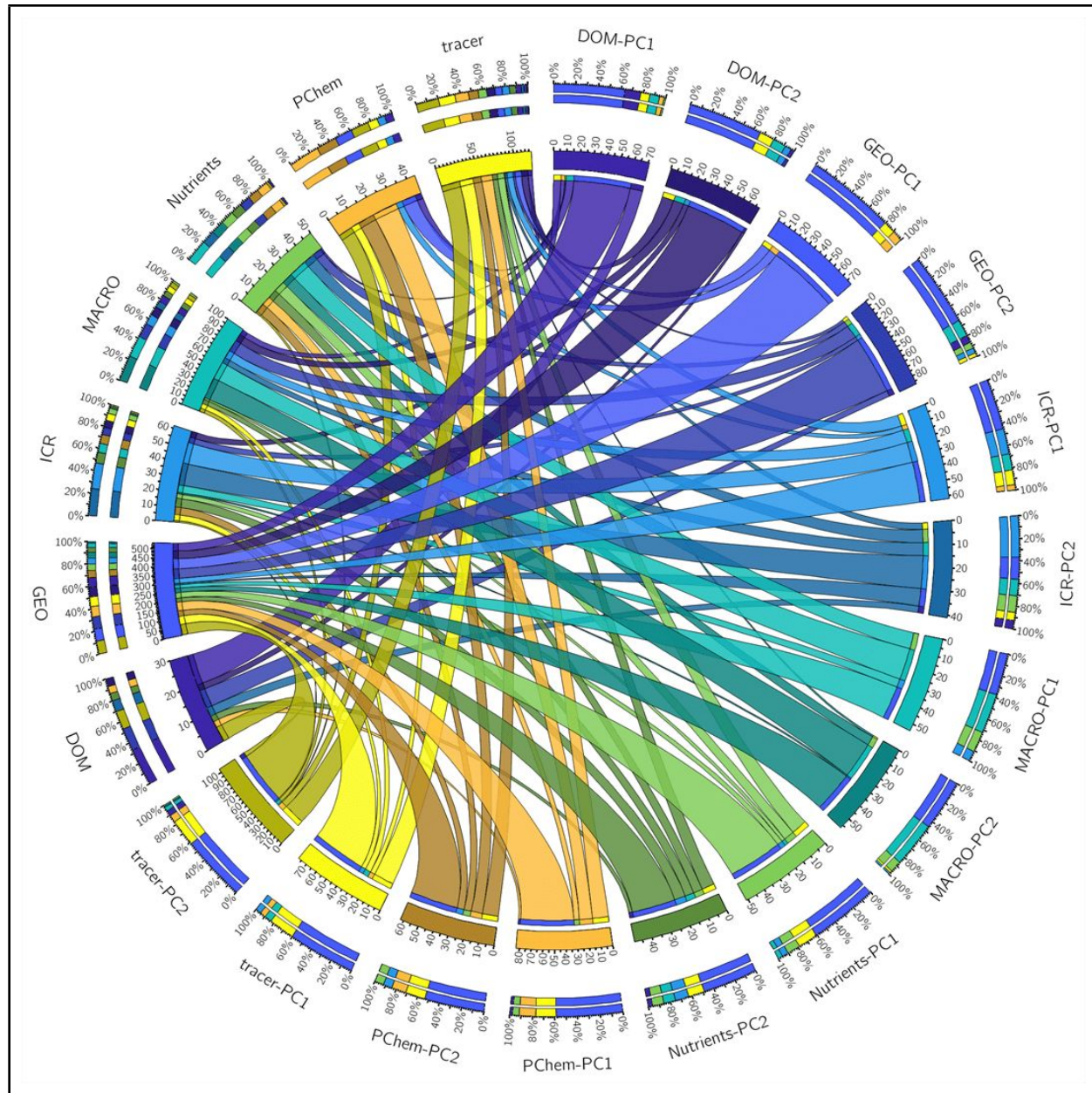


Fig. 4. Circos plot showing the one-way flow of information from the subgroup PCs (Table 1; labeled “XXX-PCY” where XXX is the subgroup and Y in the PC number) to variables predicted by the suite of SVMRs described in Section 3.3.2. Plot layout and interpretation is identical to that described for Fig. 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

1
2
3 Our literature search identified 4,075 combinations of variables that have been studied pairwise
4 in the literature (of 12,246 possible combinations). The pairwise literature search returned a total
5 of 2,731,694 results. The number of studies identified for any given pair of variables was highly
6 skewed: 50% of published studies included the 18 most commonly studied pairs of variables
7 (Ward, 2021), while the number of studies of any given pair of variables ranged from 1 to
8 270,015 (mean 670, median 14). These findings indicate a bias toward co-observation and
9 reporting of a limited number of pairwise studies, consistent with a prior study that manually
10 reviewed search results (Ward, 2015). We also found the existing literature is more focused on
11 in-group relationships (57.2% of pairwise results) compared to between-group relationships
12 (42.8% of pairwise results). In contrast, our SVMR approach identified a total of 672 pairwise
13 relationships, of which 68.8% are between-group. Notably, about 84% or 564 variable pairs do
14 not appear to have been reported previously (i.e., our systematic literature search did not return
15 any manuscripts containing information on both variables). The remaining 16% (108
16 relationships) have been previously reported in the literature (Fig. 5). The 108 relationships
17 found in both the literature and in our data analysis only represent about 2.6% of all previously-
18 reported relationships, but these relationships are included in more than 16% of all published
19 studies, indicating that prior studies have focused primarily on a relatively small number of
20 relationships. On the basis of within- and between-group frequency, the literature is broadly not
21 reflective of our findings, with the SVMR identifying higher frequencies of between-group
22 relationships that are present in the literature (Table S3). Finally, we note that the lack of a
23 relationship in the SVMR does not necessarily indicate that some relationship may be possible,
24 just as the presence of a statistical relationship does not necessarily indicate a causal relationship.
25 Some meaningful relationships could have been omitted due to signal-to-noise ratios, lagged
26 correlation between variables, or because a highly correlated variable was already selected. This
27 may explain why some well-studied relationships were not apparent in our analysis (Fig. 5).
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

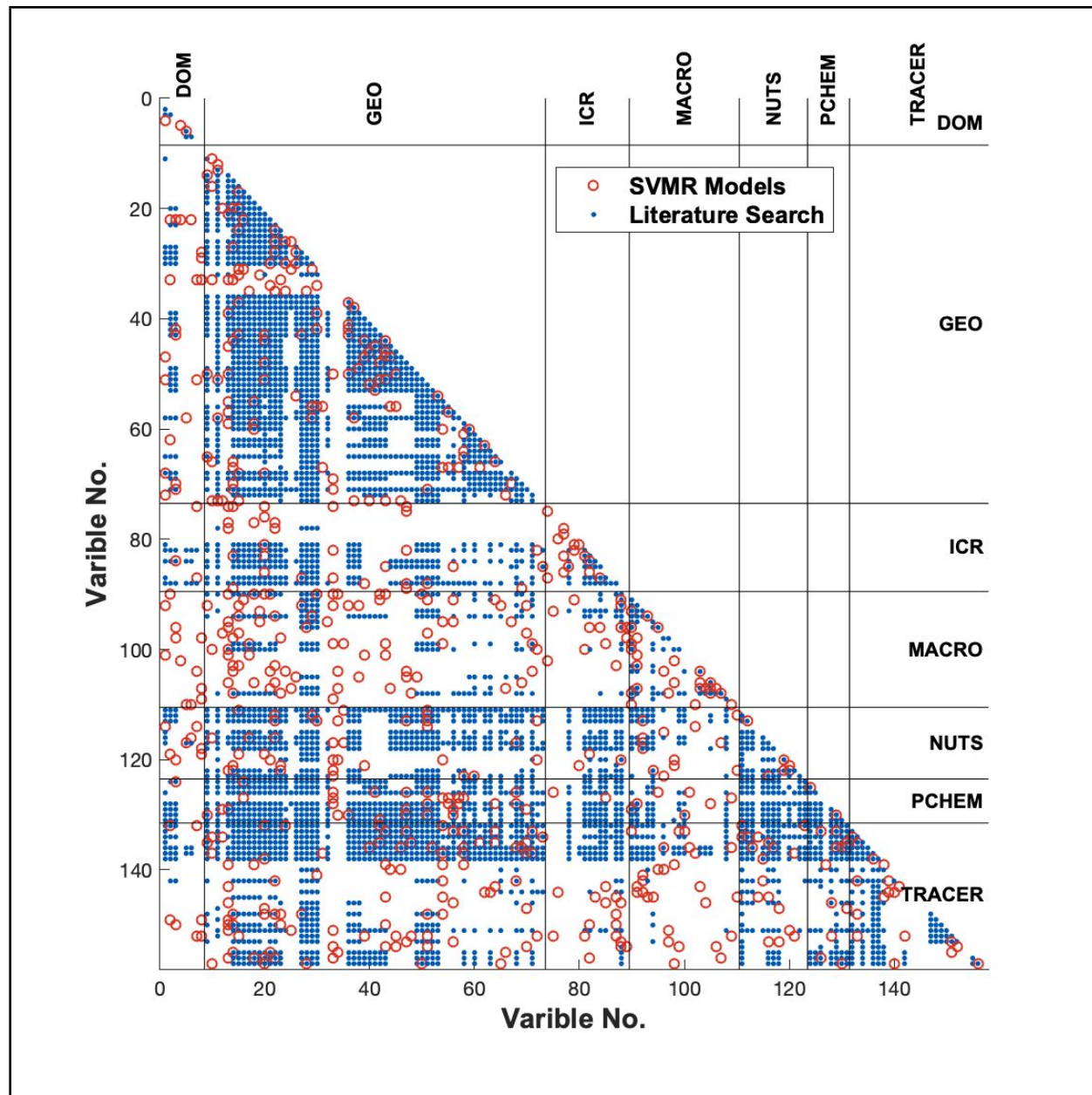


Fig. 5. Scatterplot showing pairwise study in the literature (blue dots) and identification of a relationship in our SVMR approach (red circles) for all variable pairs. Variable numbers correspond to the order variables are listed in Table S1.

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

1
2
3 27% of the variables included by SVMRs although they make up 36% of the total variable set).
4 This means the predictions of spatially structured variables were not dominated by structure from
5 a small number of structured variables. Further, a majority of variables observed (about 64%)
6 were not themselves spatially structured, and five of the seven subgroups (PCHEM, GEO,
7 NUTS, ICR, TRACER) resulted in at least one PC that was not spatially structured. These results
8 indicate that spatial structure is not ubiquitous in the river corridor. Instead, some variables
9 represent local ‘noise’ on the network-scale ‘signal’ (i.e., systematic variation in physical,
10 chemical, and biological processes from headwaters to large rivers; Vannote et al. 1980). This
11 heterogeneity is either independent from large-scale system structure (i.e., controlled by local
12 process interactions that are neither controlled by nor influence larger scale patterns) or simply
13 have sufficiently high variability to obscure larger-scale trends. Such localized ‘noise’ may also
14 reflect processes whose importance is localized in space or time, but do not recognizably follow
15 a larger spatial structure.
16
17
18
19
20
21
22
23
24
25
26

27 Individual variables reflect complex interactions that can either lead to the emergence of spatial
28 structure or overwhelm the underlying spatial structure associated with more basic variables like
29 slope and elevation. We found six variables that were spatially structured but had strong
30 relationships (SVMRs) that only included unstructured variables. In these cases, spatial structure
31 emerged or was generated by the interaction of variables that did not themselves have spatial
32 structure. Conversely, 60 of the SVMRs for unstructured variables included at least one spatially
33 structured variable (38 selected 1, 14 selected 2, and 8 selected 3 spatially structured variables).
34 This pattern suggests that spatial structure does not necessarily propagate from one variable to
35 another, indicating “signal shredding” in the river corridor (Jerolmack & Paola, 2010), where
36 information is erased by interactions between variables. While such behavior has only been
37 confirmed previously for sediment transport, our findings indicate that localized feedbacks can
38 generally overwhelm underlying spatial structure within the river corridor. This suggests that
39 sufficiently large perturbations will have system-wide impacts (e.g., large fires, floods), but
40 internal dynamics may overwhelm large-scale patterns under normal circumstances.
41 Consequently, studies of river corridors must consider local-scale interactions (i.e., internal
42 dynamics), large-scale drivers (i.e., external forcing), and the temporal context (i.e., historical
43 contingencies) if we are to account for the feedbacks and interactions in the river corridor.
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

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

1
2
3 observations (Payn et al., 2009; Valett et al., 1996) and models (Cardenas, 2008; Frissell et al.,
4 1986; Wondzell & Gooseff, 2014; Wörman et al., 2007).

5
6
7
8 Ward et al.'s (2019) observation of monotonic trends between most hydrologic exchange metrics
9 and discharge - which they describe as a proxy for network position - agree with our finding of
10 spatial structure in several variables describing geomorphic setting (including hydraulic
11 conductivity, valley slope, valley width, sinuosity), river flow (velocity, discharge), and solute
12 transport metrics (e.g., median travel time, skewness). We did not find spatial structure for other
13 metrics of exchange where Ward et al. did, including the coefficient of variation, holdback, and
14 channel water balance. Further, many of the relationships identified by Ward et al. have low
15 explanatory power as evidenced by low r^2 values, indicating that hydrologic exchange cannot be
16 described by a single explanatory variable. In contrast, the multivariate and nonlinear responses
17 encoded in the SVMs better explain the patterns in river corridor exchange observed in the
18 Andrews watersheds.
19
20
21
22
23
24
25
26
27
28

29 **4.2.2 Microbial Community Assembly**

30 Interactions along the river corridor can not only 'shred' or erase information (*sensu* Jerolmack
31 & Paola, 2010), but can also generate new information and patterns. For example, prior work at
32 the H.J. Andrews Experimental Forest spanning headwaters through 5th order rivers (Wisnoski
33 and Lennon, 2021) showed that microbial assemblages in headwater streams habitat-dependent,
34 while the microbial community became more homogeneous with distance downstream.
35 Additionally, the same study found taxonomic β -diversity was explained by an axis with positive
36 loadings for elevation and dissolved organic carbon, and negative loadings for electrical
37 conductivity, pH, total nitrogen, and total phosphorus. Microbial assemblages are known to arise
38 in response to local heterogeneity in the landscape, integrating inputs and environmental
39 variables in space and time. While we did not analyze microbial assemblages explicitly here, we
40 can interpret our observations in the context of prior findings at the site (Wisnoski and Lennon,
41 2021). Our results show spatial structure in electrical conductivity and several geomorphic
42 variables that are known to vary with elevation, but no spatial structure in total dissolved
43 phosphorus, DOC, or total dissolved nitrogen. In comparison to the controls on taxonomic β -
44 diversity described by Wisnoski and Lennon (2021), we did find spatial structure in elevation, in-
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 stream nitrate+nitrate, and electrical conductivity, but not in bulk dissolved organic carbon,
4 ammonia, or total phosphorous. Thus, our findings are broadly consistent with past findings that
5 at least some of the known controls on microbial diversity are spatially structured. However, we
6 also note that not all controls were structured, but the related microbial community did retain
7 spatial organization. Consequently, studies focused at single locations along a stream may be
8 missing contextual information on controlling factors that have propagated from the catchment
9 headwaters, or misinterpreting signals that were generated within the river corridor itself.

17 **4.2.3 River Continuum Concept**

18 The River Continuum Concept (Vannote et al., 1980) -- one of the most widely recognized and
19 cited conceptual model of river corridors -- argues that Leopold's conceptual model that
20 geomorphology reflects energy equilibrium can be extended into ecosystem functions (Langbein
21 & Leopold, 1966; L B Leopold et al., 1964; Luna B. Leopold & Langbein, 1962). Vannote et al.
22 (1980) specifically proposed: (a) biological communities should evolve to optimize the use of
23 available energy (i.e., biodegradable organic matter); and (b) energy availability will vary
24 systematically from headwaters to large downstream rivers. Our PCA results on all variables are
25 broadly consistent with these hypotheses, which is to be expected at the H.J. Andrews
26 Experimental Forest was one of the key sites studied in developing and demonstrating the
27 conceptual model. We found organic matter chemistry and geological setting explained 37% of
28 the variance across the entire data set (PC1 and PC2; Table 1). With regard to biological
29 communities optimizing to use available energy in an organized fashion, we do see that available
30 energy – in this case assessed via energy availability in organic carbon (PC1 on all variables) –
31 defines one critical dimension of variation in the system. Additionally, the high proportion of
32 spatially organized variables in TRACER, GEO, and NUTS is consistent with broad concepts of
33 systematic organization along river networks. Indeed, we found spatial structure in about 36% of
34 all variables across all disciplinary subgroups, consistent with the idea that large-scale gradients
35 will drive systematic trends in both physical and biogeochemical processes. We did find spatial
36 organization in shredders which is consistent with the River Continuum Concept. Our findings
37 on the importance of organic carbon as an explanatory variable for patterns in the river corridor
38 also support Vannote et al.'s expectation of the importance of energy availability to the structure
39 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, SUVA₂₅₄) 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 Becking 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

1
2
3 will occur, and these communities are responsible for the molecular form of organic matter that
4 is transformed within and exported from a given location. This is, functionally, the River
5 Continuum Concept applied to microbial communities. We expect the role of microbial
6 community structure in defining ecosystem processes will be critical as we transition from
7 conceptual models based on bulk measurement of organic matter (e.g., DOC, EEMs) to models
8 informed by metabolomics.
9
10
11
12
13
14

15 Previously developed theories based on bulk DOC or proxies for organic matter quality must be
16 revisited, because the field of metabolomics is rapidly evolving. The limited suite of studies that
17 include both organic carbon chemistry and nutrient data (ICR and NUTS) make comparisons for
18 consistency of findings limited. It is possible that previous conclusions about carbon limitations
19 in some systems may have been biased by only considering bulk DOC or DIC instead of its
20 molecular composition, which is highly nonuniform in its ecological function. We do not expect
21 that organic matter molecular composition is entirely controlled by geologic setting (though such
22 control has been reported; e.g., Robertson et al., 2019; Cotrufo et al., 2013), but instead that in-
23 stream organic matter reflects the integration of physical, chemical, and biological processes
24 occurring upstream of the sampling location. These processes are diverse, spanning the
25 influences of terrestrial vegetation, soil-forming processes, photochemistry, organo-mineral
26 interactions, and in-stream biological production and transformation of organic molecules. Thus,
27 the core questions are to understand when, where, and how organic matter is produced,
28 transformed, and transported. We expect that understanding microbial communities and their
29 metabolism will be critical to answering these questions.
30
31
32
33
34
35
36
37
38
39
40
41
42

43 In addition, Danczak et al. (2020) proposed a conceptual framework that draws parallels between
44 organismal birth, death, and dispersal and organic matter production, transformation, and
45 transport. They argue that organic molecules are assembled into metabolomes via a combination
46 of production, transformation, and transport just as organisms are assembled into communities
47 via a combination of birth, death, and dispersal. Danczak et al. (2020) also provide an analytical
48 approach for quantifying assembly processes, including the ability to infer when transport
49 overwhelms influences of production and transformation. This approach may be fruitful in
50 linking upland dynamics to aquatic dynamics (Waring et al., 2020; Wisnoski et al., 2021),
51
52
53
54
55
56
57
58
59
60

1
2
3 linking microbial community assembly processes to organic matter assembly processes, and
4 further highlights the need for conceptual synthesis in the river corridor (Stegen et al., 2018).
5
6
7

8 Finally, metabolomics data has been used previously to inductively reveal limitations of using
9 bulk water chemistry in river corridors to understand specific biogeochemical conditions. For
10 example, there has been a recent revelation that conceptual models for denitrification in river
11 corridors were framed at a large river network scale and not capturing dynamic, small scale
12 controls of anaerobic metabolic pathways, including denitrification (e.g., Briggs et al., 2015).
13 Since this revelation, field experiments and deductive methods have revealed that denitrification
14 is in fact occurring in sediment “microzones” across a wide range of river corridor conditions
15 that was previously hidden by and assumed impossible based upon bulk water chemistry (e.g.,
16 Knapp et al., 2017; Hampton et al., 2019; Hampton et al., 2020).
17
18
19
20
21
22
23
24

25 ***4.3.2 What controls nitrogen-acquiring extracellular enzymatic activity in a nitrogen-limited*** 26 ***ecosystem?*** 27

28 Aquatic ecosystems at the H.J. Andrews have been historically considered to be nitrogen limited
29 (Sollins et al., 1981; Triska et al., 1984). Consequently, we expected that microbes would
30 generate both leucine aminopeptidase (LAP) and N- acetylglucosaminidase (NAG) to acquire
31 nitrogen and that this would be ubiquitous across the basin. Moreover, C:N:P ratios of
32 extracellular enzymatic activity (EEA) should indicate an overproduction of N-acquiring
33 enzymes as N-limited microbes allocate energy to acquiring their limiting nutrient (e.g.,
34 Sinsabaugh et al., 1997) .
35
36
37
38
39
40
41
42

43 To test this expectation, we considered two nitrogen-acquiring enzymes: LAP and NAG. LAP
44 was part of PC1 for the NUTS subgroup and was orthogonal to total organic matter in the
45 sediment, indicating little control on sediment organic matter in explaining LAP. SVMRs for
46 LAP identify several GEO variables (bedrock type, hillslope stability, and channel water
47 balance), allochthonous inputs to the river (deciduous forest, abundance of collector-gatherer
48 macroinvertebrates), and organic carbon (spectral slope and ICR ‘other molecules’). Positive
49 correlations with spectral slope and small molecules in the ICR indicate increased LAP occurs
50 where relatively small and non-aromatic carbon sources are present. Similarly, NAG was
51
52
53
54
55
56
57
58
59
60

1
2
3 predicted by bedrock type, ICR (protein abundance), and phosphorus-acquiring enzymes.
4
5 Because we do not see spatial structure in LAP, NAG, nor 11 of the 13 variables selected by
6
7 their SVMRs, we infer that there is not a spatial control on nitrogen acquiring enzymes.
8
9

10 Several studies have reported increasing EEA with nutrient availability (Hill et al., 2010;
11
12 Sinsabaugh et al. 1997; Williams et al. 2010; Williams et al. 2012), which is not consistent with
13
14 our findings (i.e., no measurement of bulk nitrogen, carbon, phosphorus, nor oxygen were
15
16 selected by SVMRs for the ICR subgroup). Instead, we find that EEA may be explained by
17
18 particular classes of organic matter – specifically smaller, less aromatic carbon molecules,
19
20 consistent with Williams et al. (2012) and Hill et al. (2010). We also hypothesize the prevalence
21
22 of GEO features selected by SVMRs but lack of spatial structure may indicate that there are
23
24 geogenic micronutrient controls on the localized enzymatic activity that have not been measured,
25
26 such as the availability of potassium, manganese, iron, and silica that weathers from local features.

27
28 Another enzymatic question that requires more deductive work is whether the entire river
29
30 corridor is N-limited. Ecoenzymatic ratios of 1:1:1 C:N:P suggest an equilibrium between
31
32 microbial biomass and detrital organic matter (Sinsabaugh et al., 2009). The ratios of C:N and
33
34 C:P acquiring enzymes in our study (GLU:LAP+NAG and GLU:AP, respectively, based on data
35
36 in Ward et al., 2019) have slopes that are statistically indistinguishable from analyses of global
37
38 datasets (Sinsabaugh and Shah, 2012), indicating EEA is produced in relative proportions to the
39
40 basic C:N:P ratios required by microbes, suggesting that the sediment microbial community may
41
42 not, in fact, be N-limited relative to the availability of other nutrients and substrates. Therefore,
43
44 while catchment-scale mass balances indicated one understanding of the system as N-limited
45
46 (e.g., Sollins et al., 1981; Triska et al., 1984), we interpret the EEA data as an indicator that the
47
48 microbial community has adapted to the available N, and that this is present across the network
49
50 (based on the lack of spatial structure).

51
52 Our analyses suggest many fruitful paths forward for interdisciplinary river corridor research.
53
54 These include, but are not limited to, the examples presented above that (a) relate molecular
55
56 characterization of carbon to EEA to investigate organic matter quality controls; (b)
57
58 comprehensively sample stream, streambed sediment, hyporheic pore water, and hyporheic
59
60

1
2
3 sediment communities for EEA to test our hypotheses that microbes are not N limited across
4 these spatial domains; and (c) use repeated measurements to assess if one spatial snapshot of the
5 network adequately captures temporally dynamic behavior (as was found in Giraldo et al., 2014).
6
7 Our findings also suggest that the concept of ecological stoichiometry and nutrient limitations
8 manifest differently across multiple scales, warranting consideration of the places, times, and
9 scales at which equilibrium or limitation should be inferred, and whether findings of limitations
10 at one scale can be directly transferred to other scales. One particularly compelling question
11 resulting from our work is whether system-wide, large-scale N-limitation indicate low N inputs
12 at all scales, internal limitations due to spatial structure or heterogeneity (e.g., localized inputs
13 from N-fixing alders), biogeochemical limitations (e.g., kinetics of organic matter breakdown),
14 or transport limitation (e.g., inaccessibility of nutrients in some locations)?
15
16
17
18
19
20
21
22
23

24 **4.4 Inductive relationships are observations around which hypotheses can be spun and** 25 **tested**

26
27 The suite of models we constructed include 672 apparent relationships, 84% of which have not
28 been previously studied based on our literature search. It is important to recognize the
29 relationships identified here are intended as future directions, not as endpoints that reflect a
30 causal or mechanistic understanding, particularly in the case of correlations that have not been
31 reported by other studies. Each relationship serves as a set of observations, the first step in the
32 scientific method. We envision the next step for each relationship being the generation of
33 hypotheses that propose mechanisms or explanations, followed by rigorous investigation with
34 deductive approaches to rule out spurious correlation and other errors. While we have now used
35 a coarse sieve to identify mathematically meaningful relationships in the data, additional study is
36 needed to test the validity of each apparent relationship.
37
38
39
40
41
42
43
44
45

46 Even without additional investigation, it is perhaps surprising that so many apparent
47 relationships identified by our inductive approach were not found in the literature search.
48 Critically, without future study of hypotheses that can explain each relationship, like the few
49 explored in Section 4.3, we cannot differentiate if the relationships are meaningful or spurious.
50 In this regard, the inductive approach has fulfilled the promise of sieving nearly 25,000 potential
51 relationships and identifying the 672 that warrant further scrutiny. While 108 of these have been
52
53
54
55
56
57
58
59
60

1
2
3 previously reported in the literature, we identify four possibilities to explain the lack of
4 consideration of the remaining 564 pairwise statistically significant couplings in prior studies,
5 and reflect on how these results can be used to advance our goal of synthetic science to yield
6 comprehensive descriptions of the structure and function of river corridors.
7
8
9

10 11 12 ***4.4.1 Disciplinary, deductive science is the predominant mode of inquiry***

13 The norms of classical research funding opportunities and publications require deductive
14 approaches, where the limited resources of time and financial support are focused on testing
15 specific, mechanistic hypotheses. Consequently, researchers tend to dedicate effort and resources
16 on a narrow suite of specific observations rather than broader datasets that may inform the
17 connections between disciplines and scales. However, this paradigm is shifting with emphasis on
18 macrosystems research (Heffernan et al., 2014), the explicit design of networks to facilitate
19 synthesis (e.g., AmeriFlux, NEON, Critical Zone Collaborative Networks), and new funding
20 initiatives. Our results show that the inherent complexity of river corridors and networks means
21 that experimental programs of limited scope will often miss important process controls. This
22 finding provides further support for our earlier recommendation that all river corridor studies
23 collect a standard set of observations for fundamental system characterization (Ward, 2015), as
24 this information is likely to be important to testing hypotheses in ways that may not be apparent
25 in the initial study design. In this context, the inductive approach we propose here is extremely
26 useful for rapidly identifying relationships spanning disciplinary boundaries that would
27 otherwise take decades of disciplinary inquiry to identify.
28
29
30
31
32
33
34
35
36
37
38
39
40

41 42 ***4.4.2 Existing data sets are incomplete and could not have uncovered relationships***

43 Our analysis relies on the most comprehensive catchment-scale observations of interacting
44 physical, chemical, and biological processes in any river corridor to-date. The dataset we
45 analyzed also builds upon extensive prior work and data from the H.J. Andrews Experimental
46 Forest. Such comprehensive datasets, particularly co-located with long term ecological research,
47 have not previously been available and require extensive interdisciplinary collaboration to
48 obtain. For example, molecular organic matter chemistry (e.g., FTIRCMS) is only recently
49 emerging as part of river corridor science (Graham et al., 2018; Stegen, Johnson, et al., 2018;
50 Zhou et al., 2019) and has not been jointly collected with the breadth of observations we
51
52
53
54
55
56
57
58
59
60

1
2
3 analyzed here. To make further progress in unraveling the complexity of river corridors, we
4 recommend combining standardized system characterization across many streams and rivers with
5 intensive study of select watersheds to generate the rich datasets needed to evaluate process
6 interconnections and scale dependencies (Stegen & Goldman, 2018). In this case, the
7 comprehensive nature of the data set explains why novel relationships were identified here: such
8 breadth of data were simply not collected in past efforts. This further demonstrates the utility of
9 inductive analysis in generating hypotheses from new datasets that can then be tested more
10 broadly. Finally, note that our own data set, while comprehensive, is far from complete in terms
11 of all variables that could be measured across all relevant spatial scales, temporal scales, and
12 process dynamics.
13
14
15
16
17
18
19
20
21

22 ***4.4.3 Relationships may be scale- or time-dependent***

23 Both the structure and function of river corridors are known to be scale-dependent (Frissell et al.,
24 1986; Rodríguez-Iturbe & Rinaldo, 1997; McCluney et al., 2014). The network scale considered
25 here is larger than many studies of river corridors (see reviews by Tank et al., 2008; Ward,
26 2015). It is possible that the relationships identified between variables here by SVMR do not
27 hold at all scales, or that the relationships are real but have not been tested over the range of
28 scales we included in our analysis. Prior studies of river structure have found that self-
29 similarities and scale dependencies generally only occur over a limited range of scales, and either
30 average out at large scales or are limited by a physical constraint (e.g., water depth, channel
31 width, valley width) (Jerolmack & Paola, 2010; Nikora & Hicks, 1997; Rodríguez-Iturbe &
32 Rinaldo, 1997). As with relationships between individual variables, scale dependencies and
33 scaling limits identified from broad data analysis must be considered as hypotheses and tested
34 using directed observations and/or simulations with competing or alternative formulations.
35 Similarly, analyses here focused on a data set collected under baseflow conditions and process
36 controls are expected to vary in response to seasonal and storm dynamics in forcing. Moreover,
37 our analysis are focused on what can be gleaned from a single snapshot in time, whereas the
38 actual characterization includes a combination of variables spanning relatively dynamic (e.g.,
39 dissolved oxygen) to relatively static (e.g., valley slope), which may cause some relationships to
40 manifest and obscure others. Future efforts to combine high temporal resolution data with spatial
41 synoptic campaigns could directly address this limitation.
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

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.

1
2
3
4
5 Our study suggests that one avenue toward progress in river corridor science, complementary to
6 common deductive approaches, is through the collection of uniform metadata and even
7 observations typical of other scientific domains as part of disciplinary studies. We demonstrate
8 here that, in the dataset we collected, out-of-group (i.e., cross-disciplinary) data were important
9 to explaining many of the disciplinary (i.e., in-group) patterns that were observed. Thus, the out-
10 of-group data not only enable synthesis, but also simultaneously improve disciplinary
11 understanding by facilitating the generation and testing of new hypotheses. While the concepts of
12 uniform metadata and common observations have been previously called for (Ward, 2015; Ward
13 & Packman, 2019), our study demonstrates the value of these data to improve prediction of
14 individual variables or functions in the river corridor. One potentially valuable path forward
15 would be comprehensive characterization of several river corridors and at multiple times of year
16 (i.e., a modern and disciplinary broader take on the work underpinning the River Continuum
17 Concept; Minshall et al., 1983) to help determine which of the relationships we putatively
18 identify here are fundamental and general, spurious, time-variable, or organized by larger
19 climatic or geologic patterns. Another useful approach would be to identify and collect a small
20 number of variables that are informative across many sub-disciplines, and organize the findings
21 into spatially and temporally comprehensive datasets (e.g., Tiegs et al., 2019; Stegen and
22 Goldman, 2018).

23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38 In this study, we have demonstrated an application of machine learning approaches to generate
39 relationships that may inspire new studies to reveal the ‘connective tissue’ linking our
40 understanding across spatiotemporal scales and disciplines. Indeed, the step of organizing raw
41 observations to develop testable hypotheses is at the core of the scientific method, and we have
42 prototyped one approach to organize observations and highlight potential relationships in the
43 data. Hypothesis generation is touted as one of the core values of field-based observation and
44 monitoring (Burt & McDonnell, 2015; Lovett et al., 2007), where observations demand
45 explanations. The inductive approach used here presents a body of putative relationships for
46 subsequent study, at least some of which are consistent with prior conceptualizations and
47 observations of river corridors (section 4.2) and emerging areas of inquiry (section 4.3). We do
48 not propose that such approaches supplant deductive science, but rather that the two approaches
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 must be coupled in river corridor science. The inductive approach provides an unbiased or naive
4 data synthesis, which has the potential to reveal patterns and relationships that would not be
5 obvious from our present, disciplinary perspectives.
6
7
8
9

10 **5. Conclusions**

11 We began with the assumption that all variables may interact with all other variables, yielding
12 nearly 25,000 potential pairwise relationships between variables. Using machine learning, we
13 rejected most of these relationships, identifying 672 apparent relationships that have explanatory
14 power in the data set, notably including 564 pairwise relationships that were not previously
15 explored in the literature. Put another way, we have generated a web of 564 new apparent
16 relationships that may reveal new couplings in the river corridor. These relationships eschew
17 disciplinary or method-specific approaches, providing ‘connective tissue’ between traditional
18 discipline-, scale-, site-, or method-dependent knowledge. Moreover, the network of
19 relationships we have identified is consistent with several past studies from the field site
20 (Vannote et al., 1980; Ward, Wondzell, et al., 2019; Wisnoski & Lennon, 2021), providing
21 confidence that at least some of these relationships are more than spurious correlations.
22
23
24
25
26
27
28
29
30
31

32 Most of the relationships we identified, including a majority of those not present in the literature,
33 include between-group flows of information. Our results show that interactions between
34 processes that are typically studied by different disciplines is critically important to explain
35 structure and function in the river corridor. This conclusion is, perhaps, unsurprising as a
36 macrosystems view would acknowledge and expect to find cross-scale and interdisciplinary
37 relationships (Heffernan et al., 2014; McCluney et al., 2014). Still, this view is seldom fully
38 captured in existing experimental designs and the resulting data sets and literature. Importantly,
39 we also demonstrated that spatial structure can be both generated through the interaction of
40 unstructured data as well as destroyed or overprinted along the network. Thus, consideration of
41 how an observed pattern may emerge or not be visible along a spatial gradient is a critically
42 important consideration prior to interpretation of data sets.
43
44
45
46
47
48
49
50
51
52

53 Building connections between existing studies requires explicitly planning for synthesis in future
54 efforts. Here, we demonstrated the value of collecting data sets that enabled synthesis within and
55
56
57
58
59
60

1
2
3 between locations, disciplines, and scales. This does not diminish the value of traditional,
4 disciplinary hypothesis testing and deductive approaches to science. Instead, common metadata
5 and even a small number of out-of-group observations may enable synthesis efforts based on
6 inductive approaches that aids in spinning new hypotheses. Ultimately, inductive approaches are
7 a useful way to generate hypotheses from existing observational datasets and advance our
8 scientific understanding.
9
10
11
12

13 **Acknowledgements.**

14 This research has been supported by the Leverhulme Trust (Where rivers, groundwater and
15 disciplines meet: a hyporheic research network), the UK Natural Environment Research Council
16 (grant no. NE/L003872/1), the European Commission, H2020 Marie Skłodowska-Curie Actions
17 (HiFreq, grant no. 734317), the U.S. Department of Energy (Pacific Northwest National Lab and
18 DE-SC0019377), the National Science Foundation (grant nos. DEB-1440409, EAR-1652293,
19 EAR-1417603, and EAR-1446328), the University of Birmingham (Institute of Advanced
20 Studies), and with resources from the home institutions of the authors. Data and facilities were
21 provided by the H. J. Andrews Experimental Forest and Long Term Ecological Research
22 program, administered cooperatively by the USDA Forest Service Pacific Northwest Research
23 Station, Oregon State University, and the Willamette National Forest. In lieu of detailed author
24 contributions, we report that this study was conceptualized approximately 10 years ago and has
25 benefited tremendously from discussions with a broad group of friends and collaborators. Work
26 on this manuscript was initiated at the slow freshwater science meeting hold in Santa Maria de
27 Palautordera (Catalonia, NE Spain). The authors of this study each made specific contributions
28 to conceptualization, data collection, analysis, and/or writing and revising the manuscript. The
29 primary data analyzed are described by Ward et al. (2019) and available in Ward (2019). Results
30 of analyses completed in this study are available in Ward (2021). The authors declare no
31 conflicts of interest. Any use of trade, firm, or product names is for descriptive purposes only
32 and does not imply endorsement by the US government. Any opinions, findings, and conclusions
33 or recommendations expressed in this material are those of the authors.
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

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>

- 1
2
3 Frissell, C. A., Liss, W. J., Warren, C. E., & Hurley, M. D. (1986). A hierarchical framework for
4 stream habitat classification: Viewing streams in a watershed context. *Environmental*
5 *Management*, 10(2), 199–214.
- 6 Giraldo, L., Palacio, C., & Aguirre, N. (2014). Temporal Variation of the Extracellular
7 Enzymatic Activity (EEA): Case of Study : Aburra-Medellín River, in the Valle de Aburra
8 in Medellin, Antioquia, Colombia. *International Journal of Environmental Protection*, 4(5),
9 58–67.
- 10
11 Graham, E. B., Crump, A. R., Kennedy, D. W., Arntzen, E., Fansler, S., Purvine, S. O., et al.
12 (2018). Multi 'omics comparison reveals metabolome biochemistry, not microbiome
13 composition or gene expression, corresponds to elevated biogeochemical function in the
14 hyporheic zone. *Science of the Total Environment*, 642, 742–753.
15 <https://doi.org/10.1016/j.scitotenv.2018.05.256>
- 16 Gregory, K. J. (2006). The human role in changing river channels. *Geomorphology*, 79(3–4),
17 172–191. <https://doi.org/10.1016/j.geomorph.2006.06.018>
- 18
19 Hampton, TB., JP Zarnetske, MA Briggs, F MahmoodPoor Dehkordy, K Singha, FD Day-Lewis,
20 JW Harvey, S Roy Chowdhury and JW Lane. (2020) Experimental shifts of hydrologic
21 residence time in a sandy urban stream sediment–water interface alter nitrate removal and
22 nitrous oxide fluxes. *Biogeochemistry* 149, 195–219. [https://doi.org/10.1007/s10533-020-](https://doi.org/10.1007/s10533-020-00674-7)
23 [00674-7](https://doi.org/10.1007/s10533-020-00674-7).
- 24
25 Hampton, TB, JP Zarnetske, MA Briggs, K Singha, JW Harvey, FD Day-Lewis, F
26 MahmoodPoor Dehkordy, and JW Lane (2019) Residence time controls the fate of nitrogen
27 in flow-through lakebed sediments. *JGR-Biogeosciences*, 124, 689– 707.
28 <https://doi.org/10.1029/2018JG004741>
- 29
30 Harman, C. J., Ward, A. S., & Ball, A. (2016). How does reach-scale stream-hyporheic transport
31 vary with discharge? Insights fromrSAS analysis of sequential tracer injections in a
32 headwater mountain stream. *Water Resources Research*, 52, 7130–7150.
33 <https://doi.org/10.1002/2016WR018832>.Received
- 34
35 Harvey, J. W., & Gooseff, M. N. (2015). River corridor science: Hydrologic exchange and
36 ecological consequences from bedforms to basins. *Water Resources Research*, 51, 6893–
37 6922. <https://doi.org/10.1002/2015WR017617>
- 38
39 Heffernan, J. B., Soranno, P. A., Angilletta, M. J., Buckley, L. B., Gruner, D. S., Keitt, T. H., et
40 al. (2014). Macrosystems ecology: understanding ecological patterns and processes at
41 continental scales. *Frontiers in Ecology and the Environment*, 12(1), 5–14.
42 <https://doi.org/10.1890/130017>
- 43
44 Hill, B. H., McCormick, F. H., Harvey, B. C., Johnson, S. L., Warren, M. L., & Elonen, C. M.
45 (2010). Microbial enzyme activity, nutrient uptake and nutrient limitation in forested
46 streams. *Freshwater Biology*, 55(5), 1005–1019. [https://doi.org/10.1111/J.1365-](https://doi.org/10.1111/J.1365-2427.2009.02337.X)
47 [2427.2009.02337.X](https://doi.org/10.1111/J.1365-2427.2009.02337.X)
- 48
49 Isaak, D. J., Peterson, E. E., Ver Hoef, J. M., Wenger, S. J., Falke, J. A., Torgersen, C. E., et al.
50 (2014). Applications of spatial statistical network models to stream data. *Wiley*
51 *Interdisciplinary Reviews: Water*, 1(3), 277–294. <https://doi.org/10.1002/wat2.1023>
- 52
53 Jefferson, A., Grant, G. E., & Lewis, S. L. (2004). A River Runs Underneath It: Geological
54 Control of Spring and Channel Systems and Management Implications, Cascade Range,
55 Oregon. In *Advancing the Fundamental Sciences Proceedings of the Forest Service:*
56 *Proceedings of the Forest Service National Earth Sciences Conference* (Vol. 1, pp. 18–22).
- 57
58 Jerolmack, D. J., & Paola, C. (2010). Shredding of environmental signals by sediment transport.

- 1
2
3 *Geophysical Research Letters*, 37(19), 1–5. <https://doi.org/10.1029/2010GL044638>
- 4 Knapp, J. L. A., González-Pinzón, R., Drummond, J. D., Larsen, L. G., Cirpka, O. A., & Harvey,
5 J. W. (2017). Tracer-based characterization of hyporheic exchange and benthic biolayers in
6 streams. *Water Resources Res*, 53, 1575–1594. <https://doi.org/10.1002/2016WR019393>
- 7 Krause, S., Hannah, D. M., Fleckenstein, J. H., Heppell, C. M., Kaeser, D. H., Pickup, R., et al.
8 (2011). Inter-disciplinary perspectives on processes in the hyporheic zone. *Ecohydrology*,
9 4(4), 481–499.
- 10 Krause, S., Lewandowski, J., Grimm, N. B., Hannah, D. M., Pinay, G., McDonald, K., et al.
11 (2017). Ecohydrological interfaces as hot spots of ecosystem processes. *Water Resources*
12 *Research*, 53(8), 6359–6376. <https://doi.org/10.1002/2016WR019516>
- 13 Langbein, W. B., & Leopold, L. B. (1966). *River meanders - theory of minimum variance*.
- 14 Lee-Cullin, J. A., Zarnetske, J. P., Ruhala, S. S., & Plont, S. (2018). Toward measuring
15 biogeochemistry within the stream-groundwater interface at the network scale: An initial
16 assessment of two spatial sampling strategies. *Limnology and Oceanography: Methods*,
17 16(11), 722–733. <https://doi.org/10.1002/lom3.10277>
- 18 Leopold, L. B., Wolman, M. G., & Miller, J. P. (1964). *Fluvial Processes in Geomorphology*.
19 Dover Publications.
- 20 Leopold, L. B., & Langbein, W. B. (1962). *The Concept of Entropy in Landscape Evolution*.
- 21 Li, L., Sullivan, P. L., Benettin, P., Cirpka, O. A., Bishop, K., Brantley, S. L., et al. (2021).
22 Toward catchment hydro-biogeochemical theories. *Wiley Interdisciplinary Reviews: Water*,
23 8(1), e1495. <https://doi.org/10.1002/wat2.1495>
- 24 Liébault, F., & Piégay, H. (2002). Causes of 20th century channel narrowing in mountain and
25 piedmont rivers of southeastern France. *Earth Surface Processes and Landforms*, 27(4),
26 425–444. <https://doi.org/10.1002/esp.328>
- 27 Lovett, G. M., Burns, D. A., Driscoll, C. T., Jenkins, J. C., Mitchell, M. J., Rustad, L., et al.
28 (2007). Who needs environmental monitoring? *Frontiers in Ecology and the Environment*,
29 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)
- 30 Martin, P. Y., & Turner, B. A. (1986). Grounded Theory and Organizational Research. *The*
31 *Journal of Applied Behavioral Science*, 22(2), 141–157.
32 <https://doi.org/10.1177/002188638602200207>
- 33 McCluney, K. E., Poff, N. L., Palmer, M. A., Thorp, J. H., Poole, G. C., Williams, B. S., et al.
34 (2014). Riverine macrosystems ecology: sensitivity, resistance, and resilience of whole river
35 basins with human alterations. *Frontiers in Ecology and the Environment*, 12(1), 48–58.
36 <https://doi.org/10.1890/120367>
- 37 McGuire, K. J., Torgersen, C. E., Likens, G. E., Buso, D. C., Lowe, W. H., & Bailey, S. W.
38 (2014). Network analysis reveals multiscale controls on streamwater chemistry.
39 *Proceedings of the National Academy of Sciences of the United States of America*, 111(19),
40 7030–7035. <https://doi.org/10.1073/pnas.1404820111>
- 41 Minshall, G. W., Petersen, R. C., Cummins, K. W., Bott, T. L., Sedell, J. R., Cushing, C. E., &
42 Vannote, R. L. (1983). Interbiome Comparison of Stream Ecosystem Dynamics. *Ecological*
43 *Monographs*, 53(1), 1–25. <https://doi.org/10.2307/1942585>
- 44 Nikora, V. I., & Hicks, D. M. (1997). Scaling Relationships for Sand Wave Development in
45 Unidirectional Flow. *Journal of Hydraulic Engineering*, 123(12), 1152–1156.
46 [https://doi.org/10.1061/\(asce\)0733-9429\(1997\)123:12\(1152\)](https://doi.org/10.1061/(asce)0733-9429(1997)123:12(1152))
- 47 O'Malley, M.A. (2008). “Everything is everywhere: but the environment selects”: ubiquitous
48 distribution and ecological determinism in microbial biogeography. *Studies in History and*
49
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3 *Philosophy of Biological and Biomedical Sciences*, 39(3), 314–325.
4 <https://doi.org/10.1016/J.SHPSC.2008.06.005>
5 Payn, R. A., Gooseff, M. N., McGlynn, B. L., Bencala, K. E., & Wondzell, S. M. (2009).
6 Channel water balance and exchange with subsurface flow along a mountain headwater
7 stream in Montana, United States. *Water Resources Research*, 45.
8 Pinay, G., Peiffer, S., De Dreuzy, J. R., Krause, S., Hannah, D. M., Fleckenstein, J. H., et al.
9 (2015). Upscaling Nitrogen Removal Capacity from Local Hotspots to Low Stream Orders'
10 Drainage Basins. *Ecosystems*, 18(6), 1101–1120. [https://doi.org/10.1007/s10021-015-9878-](https://doi.org/10.1007/s10021-015-9878-5)
11 [5](https://doi.org/10.1007/s10021-015-9878-5)
12 Pringle, C. M., Naiman, R. J., Bretschko, G., Karr, J. R., Oswood, M. W., Webster, J. R., et al.
13 (1988). Patch Dynamics in Lotic Systems: The Stream as a Mosaic. *Journal of the North*
14 *American Benthological Society*, 7(4), 503–524. <https://doi.org/10.2307/1467303>
15 Rana, S. M. M., Scott, D. T., & Hester, E. T. (2017). Effects of in-stream structures and channel
16 flow rate variation on transient storage. *Journal of Hydrology*, 548, 157–169.
17 <https://doi.org/10.1016/j.jhydrol.2017.02.049>
18 Robertson, A. D., Paustian, K., Ogle, S., Wallenstein, M. D., Lugato, E., & Francesca Cotrufo,
19 M. (2019). Unifying soil organic matter formation and persistence frameworks: The MEMS
20 model. *Biogeosciences*, 16(6), 1225–1248. <https://doi.org/10.5194/BG-16-1225-2019>
21 Rodríguez-Iturbe, I., & Rinaldo, A. (1997). *Fractal River Basins: Chance and Self-Organization*.
22 Cambridge, UK: Cambridge University Press.
23 Santschi, P. H., Presley, B. J., Wade, T. L., Garcia-Romero, B., & Baskaran, M. (2001).
24 Historical contamination of PAHs, PCBs, DDTs, and heavy metals in Mississippi River
25 Delta, Galveston Bay and Tampa Bay sediment cores. *Marine Environmental Research*,
26 52(1), 51–79. [https://doi.org/10.1016/S0141-1136\(00\)00260-9](https://doi.org/10.1016/S0141-1136(00)00260-9)
27 Sinsabaugh, R. L., Findlay, S., Franchini, P., & Fischer, D. (1997). Enzymatic analysis of
28 riverine bacterioplankton production. *Limnology and Oceanography*, 42(1), 29–38.
29 <https://doi.org/10.4319/LO.1997.42.1.0029>
30 Sinsabaugh, R. L., Findlay, S., Franchini, P., & Fischer, D. (1997). Enzymatic analysis of
31 riverine bacterioplankton production. *Limnology and Oceanography*, 42(1), 29–38.
32 <https://doi.org/10.4319/LO.1997.42.1.0029>
33 Sinsabaugh, R. L., & Shah, J. J. F. (2012). Ecoenzymatic Stoichiometry and Ecological Theory.
34 <http://Dx.Doi.Org/10.1146/Annurev-Ecolsys-071112-124414>, 43, 313–343.
35 <https://doi.org/10.1146/ANNUREV-ECOLSYS-071112-124414>
36 Smidt, S. J., Cullin, J. A., Ward, A. S., Robinson, J., Zimmer, M. A., Lutz, L. K., & Endreny, T.
37 A. (2015). A Comparison of Hyporheic Transport at a Cross-Vane Structure and Natural
38 Riffle. *Ground Water*, 53(6), 859–871. <https://doi.org/10.1111/gwat.12288>
39 Sollins, P., Cromack, K., Corison, F. M. M., Waring, R. H., & Harr, R. D. (1981). Changes in
40 Nitrogen Cycling at an Old-Growth Douglas-fir Site After Disturbance. *Journal of*
41 *Environmental Quality*, 10(1), 37–42.
42 <https://doi.org/10.2134/JEQ1981.00472425001000010007X>
43 Stegen, J. C., & Goldman, A. E. (2018). WHONDRS: a Community Resource for Studying
44 Dynamic River Corridors. *MSystems*, 3(5), 151–169.
45 <https://doi.org/10.1128/msystems.00151-18>
46 Stegen, J. C., Bottos, E. M., & Jansson, J. K. (2018). A unified conceptual framework for
47 prediction and control of microbiomes. *Current Opinion in Microbiology*, 44(July), 20–27.
48 <https://doi.org/10.1016/j.mib.2018.06.002>
49
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3 Stegen, J. C., Johnson, T., Fredrickson, J. K., Wilkins, M. J., Konopka, A. E., Nelson, W. C., et
4 al. (2018). Influences of organic carbon speciation on hyporheic corridor biogeochemistry
5 and microbial ecology. *Nature Communications*, 9(1), 1–11.
6 <https://doi.org/10.1038/s41467-018-03572-7>
7
8 Strauss, A., & Corbin, J. (1994). Grounded theory methodology: An overview. In N. Denzin &
9 Y. Lincoln (Eds.), *Handbook of qualitative research* (pp. 273–285). Sage Publications, Inc.
10 Swanson, F. J., & James, M. E. (1975). *Geology and geomorphology of the H.J. Andrews*
11 *Experimental Forest, western Cascades, Oregon*. Portland, OR.
12 Swanson, F. J., & Jones, J. A. (2002). Geomorphology and hydrology of the H.J. Andrews
13 Experimental Forest, Blue River, Oregon. In *Field guide to geologic processes in Cascadia*.
14 Tank, J. L., Rosi-Marshall, E. J., Baker, M. A., & Hall, R. O. (2008). Are rivers just big streams?
15 A pulse method to quantify nitrogen demand in a large river. *Ecology*, 89(10), 2935–2945.
16 Tiegs, S. D., Costello, D. M., Isken, M. W., Woodward, G., McIntyre, P. B., Gessner, M. O., et
17 al. (2019). Global patterns and drivers of ecosystem functioning in rivers and riparian zones.
18 *Science Advances*, 5(1), eaav0486. <https://doi.org/10.1126/SCIADV.AAV0486>
19 Triska, F. J., Sedell, J. R., Cromack, K., Gregory, S. V., & McCorison, F. M. (1984). Nitrogen
20 Budget for a Small Coniferous Forest Stream. *Ecological Monographs*, 54(1), 119–140.
21 <https://doi.org/10.2307/1942458>
22
23 Turnbull, L., Hütt, M. T., Ioannides, A. A., Kininmonth, S., Poepl, R., Tockner, K., et al. (2018,
24 December 1). Connectivity and complex systems: learning from a multi-disciplinary
25 perspective. *Applied Network Science*. Springer. <https://doi.org/10.1007/s41109-018-0067-2>
26 Valett, H. M., Morrice, J. A., Dahm, C. N., & Campana, M. E. (1996). Parent lithology, surface-
27 groundwater exchange, and nitrate retention in headwater streams. *Limnology and*
28 *Oceanography*, 333–345.
29
30 Vannote, R. L., Minshall, G. W., Cummins, K. W., Sedell, J. R., & Cushing, C. E. (1980). The
31 River Continuum Concept. *Canadian Journal of Fisheries and Aquatic Sciences*, 37, 130–
32 137.
33
34 Ver Hoef, J. M., Peterson, E., & Theobald, D. (2006). Spatial statistical models that use flow and
35 stream distance. *Environmental and Ecological Statistics*, 13(4), 449–464.
36 <https://doi.org/10.1007/s10651-006-0022-8>
37
38 Walling, D. E., & Fang, D. (2003). Recent trends in the suspended sediment loads of the world's
39 rivers. *Global and Planetary Change*, 39(1–2), 111–126. [https://doi.org/10.1016/S0921-](https://doi.org/10.1016/S0921-8181(03)00020-1)
40 [8181\(03\)00020-1](https://doi.org/10.1016/S0921-8181(03)00020-1)
41
42 Wallis, I., Prommer, H., Berg, M., Siade, A. J., Sun, J., & Kipfer, R. (2020). The river–
43 groundwater interface as a hotspot for arsenic release. *Nature Geoscience*, 13(4), 288–295.
44 <https://doi.org/10.1038/s41561-020-0557-6>
45
46 Ward, A. S. (2015). The evolution and state of interdisciplinary hyporheic research. *Wiley*
47 *Interdisciplinary Reviews: Water*, 3(1), 83–103. <https://doi.org/10.1002/wat2.1120>
48
49 Ward, A. S. (2019). ESSD, 2019 - Data Collection. <https://doi.org/10.5194/essd-11-1-2019>
50
51 Ward, A. S. (2021). Supporting data for Ward et al., (In Review) Advancing river corridor
52 science beyond disciplinary boundaries with an inductive approach to hypothesis
53 generation, HydroShare, Accessed 6-May-2021.
54 <http://www.hydroshare.org/resource/de6d92d314354ea6819157818669fc59>
55
56 Ward, A. S., & Packman, A. I. (2019). Advancing our predictive understanding of river corridor
57 exchange. *Wiley Interdisciplinary Reviews: Water*, 6(1), e1327.
58 <https://doi.org/10.1002/wat2.1327>
59
60

- 1
2
3 Ward, A. S., Zarnetske, J. P., Baranov, V., Blaen, P. J., Brekenfeld, N., Chu, R., et al. (2019).
4 Co-located contemporaneous mapping of morphological, hydrological, chemical, and
5 biological conditions in a 5th-order mountain stream network, Oregon, USA. *Earth System*
6 *Science Data*, 11(4). <https://doi.org/10.5194/essd-11-1567-2019>
7
8 Ward, A. S., Wondzell, S. M., Schmadel, N. M., Herzog, S., Zarnetske, J. P., Baranov, V., et al.
9 (2019). Spatial and temporal variation in river corridor exchange across a 5th order
10 mountain stream network. *Hydrology and Earth System Sciences Discussions*, (April), 1–
11 39. <https://doi.org/10.5194/hess-2019-108>
12
13 Waring, B. G., Sulman, B. N., Reed, S., Smith, A. P., Averill, C., Creamer, C. A., et al. (2020).
14 From pools to flow: The PROMISE framework for new insights on soil carbon cycling in a
15 changing world. *Global Change Biology*, 26(12), 6631–6643.
16 <https://doi.org/10.1111/GCB.15365>
17
18 Williams, C. J., Scott, A. B., Wilson, H. F., & Xenopoulos, M. A. (2011). Effects of land use on
19 water column bacterial activity and enzyme stoichiometry in stream ecosystems. *Aquatic*
20 *Sciences 2011* 74:3, 74(3), 483–494. <https://doi.org/10.1007/S00027-011-0242-3>
21
22 Williams, C. J., Yamashita, Y., Wilson, H. F., Jaffé, R., & Xenopoulos, M. A. (2010).
23 Unraveling the role of land use and microbial activity in shaping dissolved organic matter
24 characteristics in stream ecosystems. *Limnology and Oceanography*, 55(3), 1159–1171.
25 <https://doi.org/10.4319/LO.2010.55.3.1159>
26
27 Wisnoski, N. I., & Lennon, J. T. (2021). Microbial community assembly in a multi-layer
28 dendritic metacommunity. *Oecologia*, 195(1), 13–24. [https://doi.org/10.1007/s00442-020-](https://doi.org/10.1007/s00442-020-04767-w)
29 [04767-w](https://doi.org/10.1007/s00442-020-04767-w)
30
31 Wisnoski, N. I., Muscarella, M. E., Larsen, M. L., Peralta, A. L., & Lennon, J. T. (2020).
32 Metabolic insight into bacterial community assembly across ecosystem boundaries.
33 *Ecology*, 101(4), e02968. <https://doi.org/10.1002/ECY.2968>
34
35 Wit, R. De, & Bouvier, T. (2006). ‘Everything is everywhere, but, the environment selects’; what
36 did Baas Becking and Beijerinck really say? *Environmental Microbiology*, 8(4), 755–758.
37 <https://doi.org/10.1111/J.1462-2920.2006.01017.X>
38
39 Wohl, E. (2005). Compromised Rivers: Understanding Historical Human Impacts on Rivers in
40 the Context of Restoration. *Ecology and Society*, 10(2), 2.
41
42 Wondzell, S. M., & Gooseff, M. N. (2014). Geomorphic Controls on Hyporheic Exchange
43 Across Scales: Watersheds to Particles. In J. Schroder & E. Wohl (Eds.), *Treatise on*
44 *Geomorphology* (Vol. 9, pp. 203–218). San Diego, CA: Academic Press.
45
46 Wood, P. J., Boulton, A. J., Little, S., & Stubbington, R. (2010). Is the hyporheic zone a
47 refugium for aquatic macroinvertebrates during severe low flow conditions? *Fundamental*
48 *and Applied Limnology / Archiv Für Hydrobiologie*, 176(4), 377–390.
49 <https://doi.org/10.1127/1863-9135/2010/0176-0377>
50
51 Wörman, A., Packman, A. I., Marklund, L., Harvey, J. W., & Stone, S. H. (2007). Fractal
52 topography and subsurface water flows from fluvial bedforms to the continental shield.
53 *Geophysical Research Letters*, 34(7), 1–5. <https://doi.org/10.1029/2007GL029426>
54
55 Wu, L., Singh, T., Gomez-Velez, J., Nützmann, G., Wörman, A., Krause, S., & Lewandowski, J.
56 (2018). Impact of Dynamically Changing Discharge on Hyporheic Exchange Processes
57 Under Gaining and Losing Groundwater Conditions. *Water Resources Research*, 54(12),
58 10,076–10,093. <https://doi.org/10.1029/2018WR023185>
59
60 Yoder, L., Ward, A. S., Spak, S., & Dalrymple, K. (2020). Local Government Perspectives on
Collaborative Governance: A Comparative Analysis of Iowa’s Watershed Management

1
2
3 Authorities. *Policy Studies Journal*. <https://doi.org/10.1111/psj.12389>
4 Zhou, C., Liu, Y., Liu, C., Liu, Y., & Tfaily, M. M. (2019). Compositional changes of dissolved
5 organic carbon during its dynamic desorption from hyporheic zone sediments. *Science of*
6 *the Total Environment*, 658, 16–23. <https://doi.org/10.1016/j.scitotenv.2018.12.189>
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

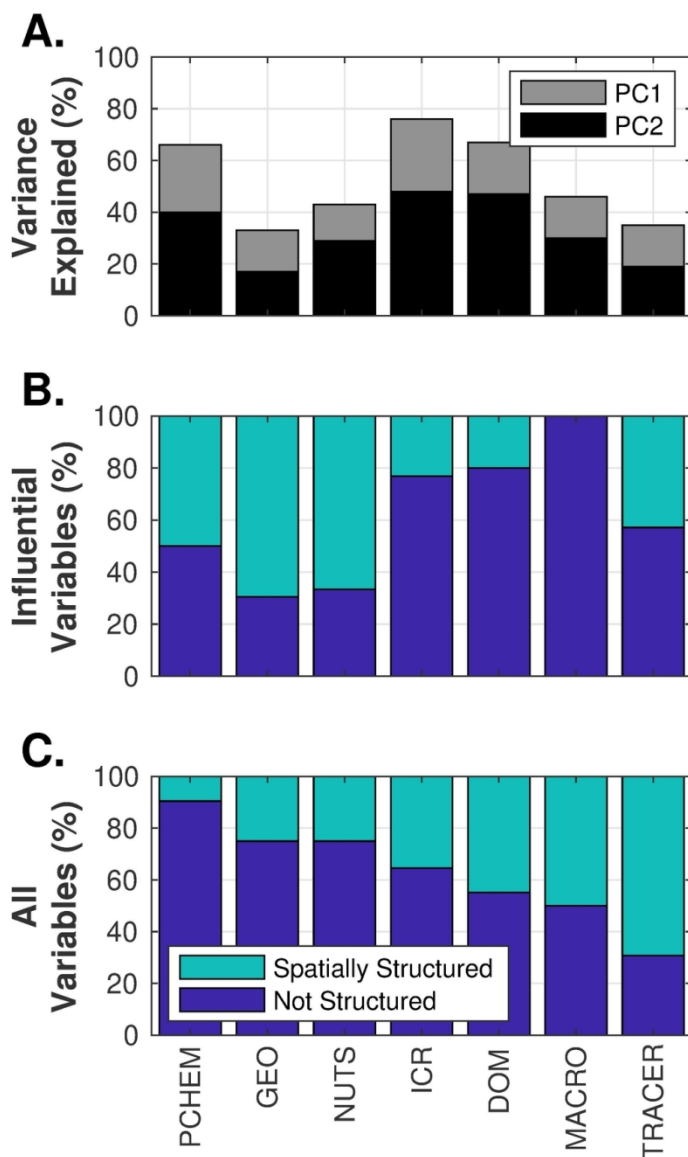


Fig. 1. (A) Variance in the Andrews river corridor data set explained by PC1 and PC2 for each expert subgroup. (B) Percentage of influential variables (i.e., the variables included in the first two PCs) that do and do not have spatial structure. (C) Percentage of all variables within each subgroup that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

94x149mm (300 x 300 DPI)

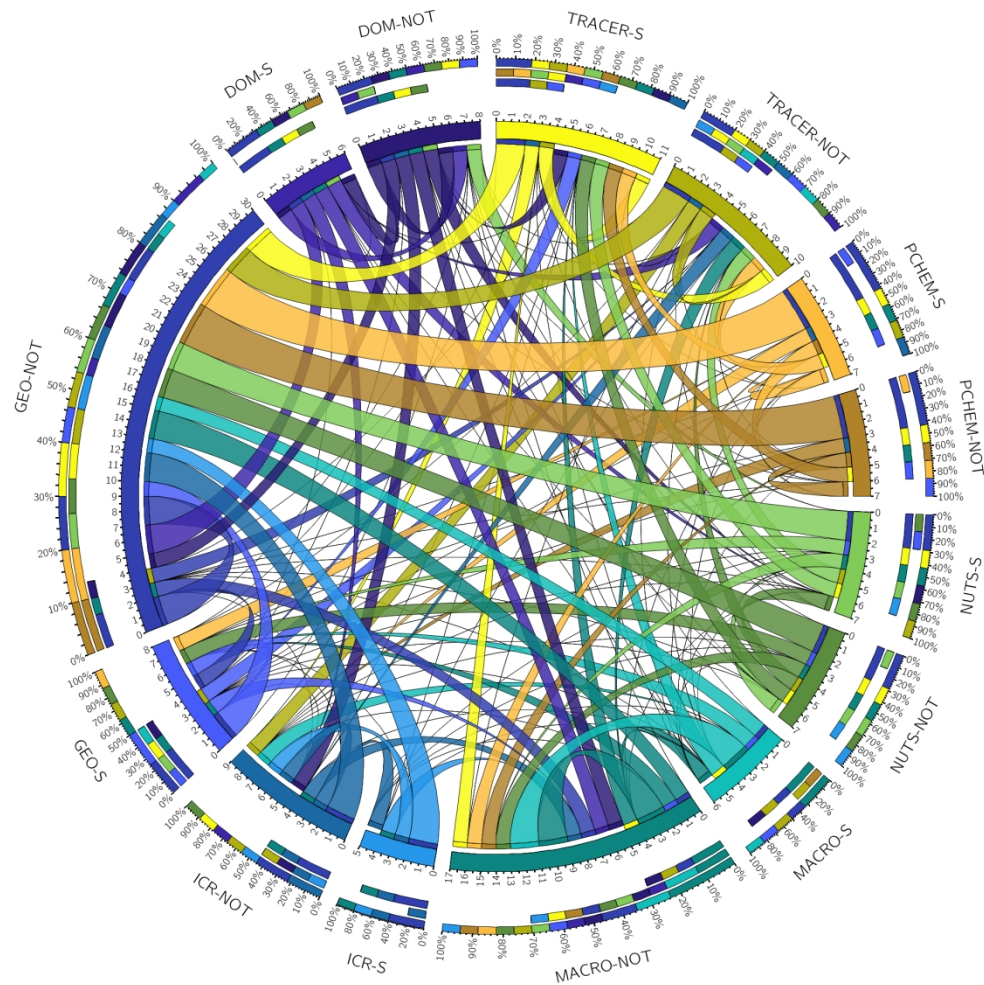


Fig. 2. Information flow within and among subgroups of variables commonly used as measures of river corridor dynamics based on the suite of SVMRs constructed for each variable (Section 3.3.1). The variables included in the 7 subgroups (PCHEM = physical chemistry; GEO = geologic setting; NUTS = nutrients; 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.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

- 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.

705x705mm (72 x 72 DPI)

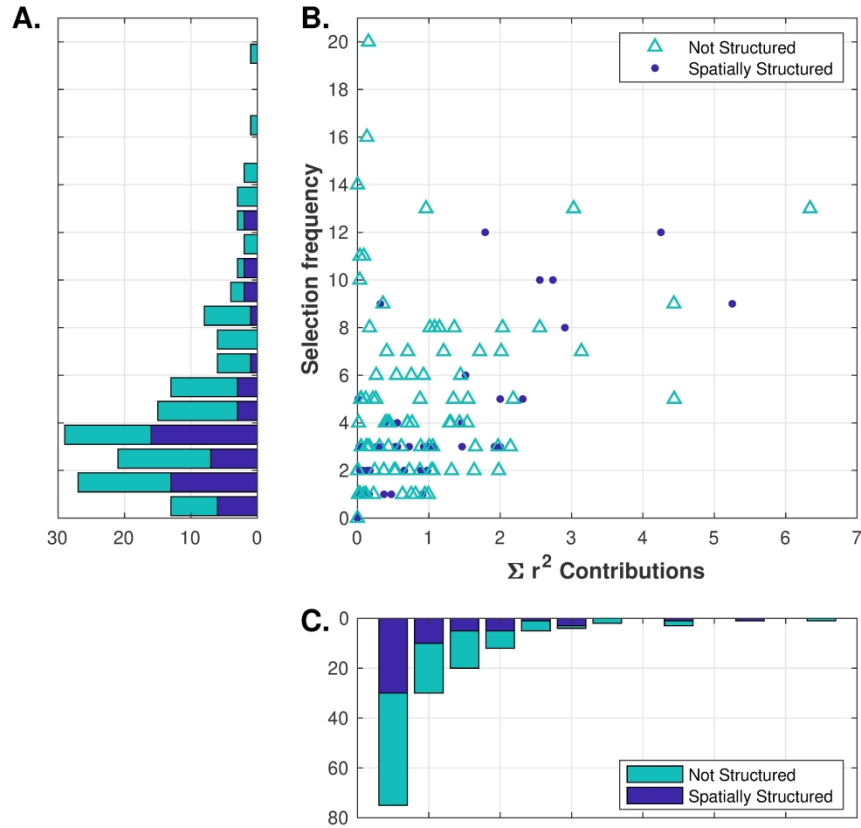


Fig. 3. Distributions of variable selection frequency and contributions of each variable to improvement in r^2 for all 157 SVRMs constructed on all variables. (A) Histogram of the frequency with which individual variables were selected. (B) Variable selection frequency vs. total improvement in r^2 . (C) Histogram of contributions of variables to r^2 .

190x190mm (600 x 600 DPI)

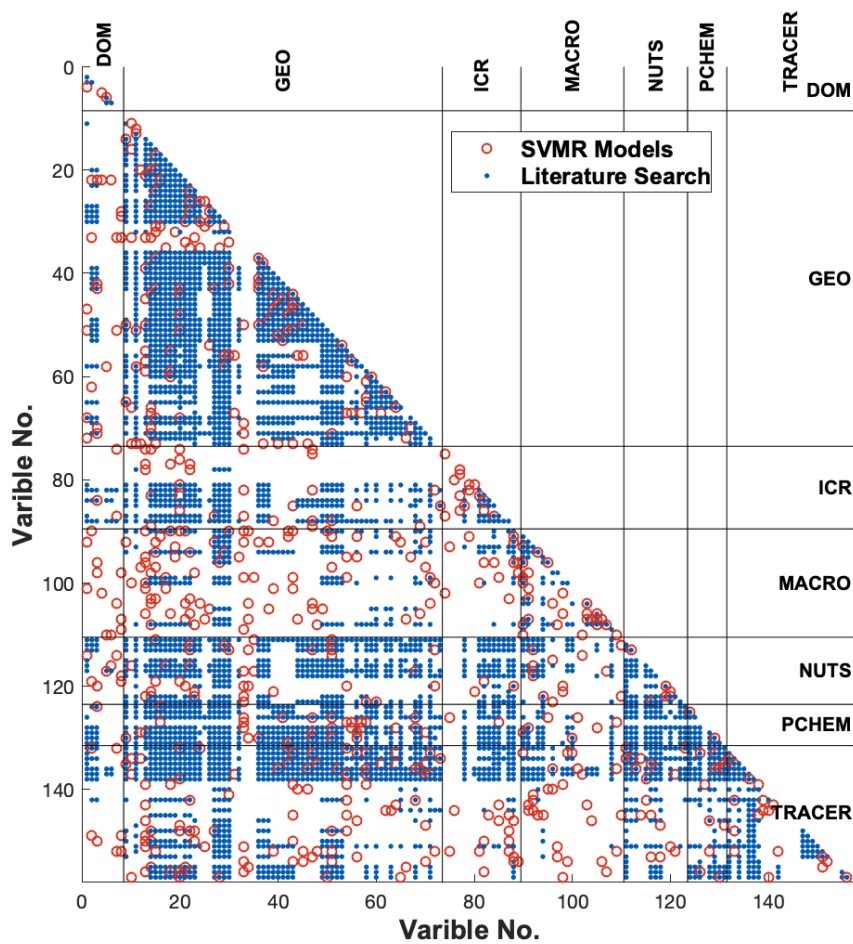


Fig. 4. Circos plot showing the one-way flow of information from the subgroup PCs (Table 1; labeled “XXX-PCY” where XXX is the subgroup and Y in the PC number) to variables predicted by the suite of SVMRs described in Section 3.3.2. Plot layout and interpretation is identical to that described for Fig. 2.

977x977mm (28 x 28 DPI)

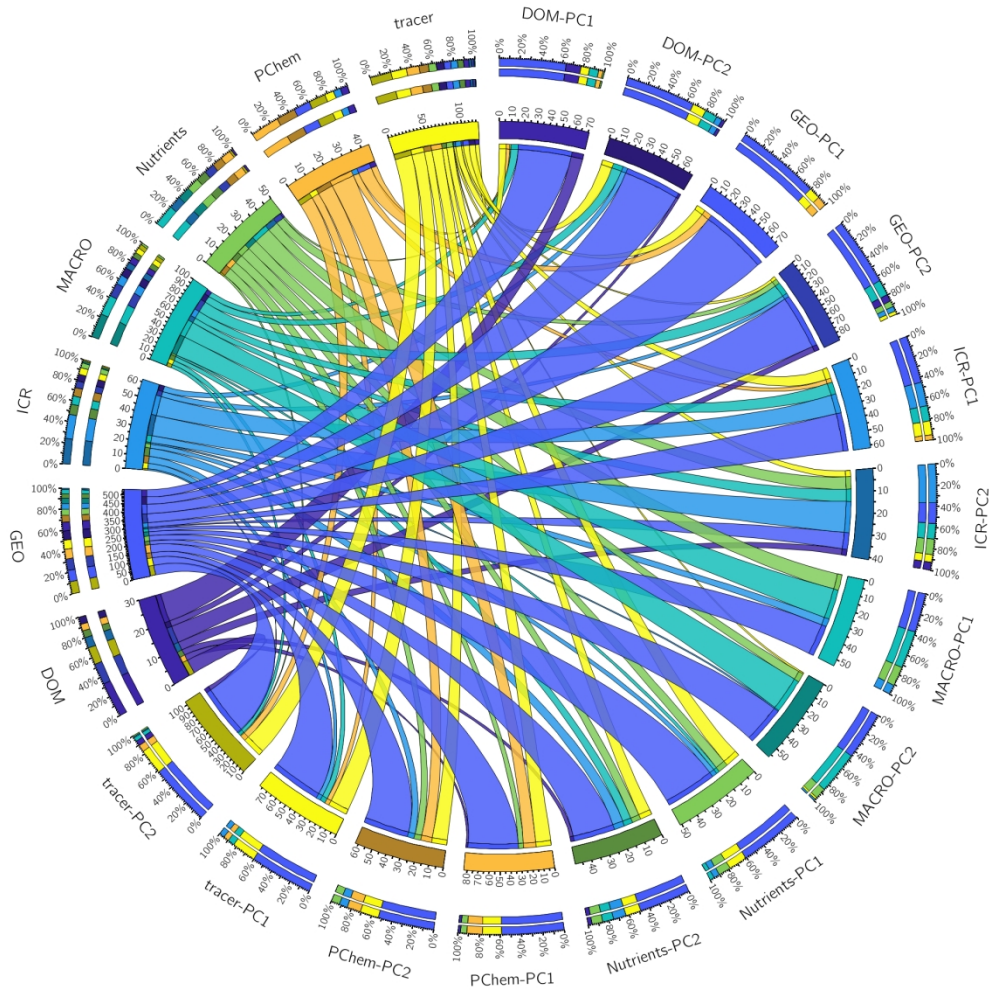


Fig. 5. Scatterplot showing pairwise study in the literature (blue dots) and identification of a relationship in our SVMR approach (red circles) for all variable pairs. Variable numbers correspond to the order variables are listed in Table S1.

705x705mm (72 x 72 DPI)

PCA on all variables

	PC1			PC2		
	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

PCA on subgroups

	PC1			PC2		
	Variance explained (%)	Positive loadings	Negative loading	Variance explained (%)	Positive loadings	Negative loading
Physical Chemistry (PCEM)	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	16	Abundance of collector-gatherers	Abundance of shredders, Abundance of small body size
Stream Solute Tracer (TRACER)	19 *	—	short term storage (holdback, skewness, CV)	16	Dispersion, Fraction of mass in A/D, velocity, upstream and downstream discharge	—

* Indicates the PC is spatially structured

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

240x172mm (200 x 200 DPI)

1
2
3 **Advancing river corridor science beyond disciplinary boundaries with an inductive**
4 **approach to catalyze hypothesis generation**
5

6 *Submitted for publication in Hydrological Processes – Special issue – Data Science Applications*
7 *in Hydrology*
8
9

10 **Authors:**

11 Adam S. Ward¹, Aaron Packman², Susana Bernal³, Nicolai Brekenfeld⁴, Jen Drummond⁴, Emily
12 Graham⁵, David M. Hannah⁴, Megan Klaar⁶, Stefan Krause⁴, Marie Kurz⁷, Angang Li², Anna
13 Lupon³, Feng Mao⁸, M. Eugènia Martí Roca³, Valerie Ouellet⁴, Todd V. Royer¹, James C.
14 Stegen⁵, Jay P. Zarnetske⁹
15

16
17 ¹ O'Neill School of Public and Environmental Affairs, Indiana University, Bloomington,
18 Indiana, USA

19 ² Department of Civil and Environmental Engineering, Northwestern University, Evanston,
20 Illinois, USA

21 ³ Integrative Freshwater Ecology Group, Centre for Advanced Studies of Blanes (CEAB-
22 CSIC), Blanes, Spain

23 ⁴ School of Geography, Earth & Environmental Sciences, University of Birmingham,
24 Edgbaston, Birmingham, B15 2TT, UK

25 ⁵ Earth and Biological Sciences Directorate, Pacific Northwest National Laboratory, Richland,
26 Washington, USA

27 ⁶ School of Geography, School of Earth and Environment, University of Leeds, Woodhouse,
28 Leeds LS2 9JT, United Kingdom

29 ⁷ The Academy of Natural Sciences of Drexel University, Philadelphia, Pennsylvania, USA

30 ⁸ School of Earth and Environmental Sciences, Cardiff University, Building, Park Place,
31 Cardiff, CF10 3AT, United Kingdom

32 ⁹ Department of Earth and Environmental Sciences, Michigan State University, East Lansing,
33 Michigan, USA
34
35
36

37
38 **Corresponding author:**

39 Adam S. Ward
40 O'Neill School of Public and Environmental Affairs
41 Indiana University
42 418 MSB-II
43 Bloomington, IN 47405
44

45 Email: adamward@indiana.edu

46 Phone: 812-865-4820
47
48

49 **Running head:** Inductive hypothesis generation using data science
50

51 **Key words:** river corridor, stream corridor, machine learning, inductive, scientific method
52
53
54
55
56
57
58

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 5th 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 for hypothesis generation, to sieving complex data sets and identify for 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.

This structured approach provides a means to unify the fragmented knowledge gained by traditional modes of inquiry.

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 ~~lack the ability~~have yet to connect-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 ~~process dynamics~~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

1
2
3 across disciplines. One way to address these needs is to expand beyond the traditional, deductive
4 approach to science, which bases measurements on a highly targeted set of causal mechanisms to
5 be tested at a limited range of locations and scales. With the emergence of new experimental and
6 data science techniques, the time has come to expand existing conceptual models for river
7 corridors via approaches that generate more integrative knowledge commensurate with the
8 reality of of river corridors as complex dynamic systems. We posit that unified understanding
9 must be derived from a combination of *deductive* science and *inductive* approaches that identify
10 process interactions and couplings that emerge from the data themselves. We suggest that river
11 corridor science can benefit from inductive approaches that generate hypotheses and eventually
12 theories from empirical studies, an approach successfully applied in other disciplines ~~Complex~~
13 ~~Systems and Grounded Theory approaches that have proven useful in understanding many other~~
14 ~~problems that involve complex multiscale dynamics~~ (Martin & Turner, 1986; Strauss & Corbin,
15 1994; e.g., Turnbull et al., 2018).

16
17
18
19
20
21
22
23
24
25
26
27 A unifying framework is required to organize and synthesize our understanding of river corridors
28 and advance scientific understanding of the drivers and controls of their functioning. Stegen et al.
29 (2018) propose one such model for microbial ecology, where the resultant ecosystem functions
30 and services are explained by the relationships linking internal dynamics, external forcing, and
31 historical contingencies. The principles of Stegen et al.'s conceptual framework are similar to
32 other existing conceptualizations of river corridors that have been developed by other disciplines.
33 First, external forcing describes the role of factors extrinsic to the river corridor that shape its
34 structure and function. For river corridors, this primarily means the larger spatial scale and
35 longer temporal scale elements that are functionally decoupled (e.g., static or slowly-varying)
36 relative to a process of interest. Studies with data collection spanning gradients in land use,
37 geologic setting, climate, network position, or other factors that are considered to be extrinsic
38 typically use geospatial and statistical approaches to describe patterns and trends (e.g., McGuire
39 et al., 2014), while variation around spatially structured trends is often interpreted as random
40 noise from structural heterogeneity and/or unstudied, smaller-scale processes (Abbott et al.,
41 2018). Next, internal dynamics are the interacting processes within the river corridor that give
42 rise to observed functions of interest at a given location. Conceptual models based on this
43 approach to river corridor science include hot spots and hot moments (Krause et al., 2011, 2017;
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Wallis et al., 2020), control points (Bernhardt et al., 2017), and patch dynamics (Pringle et al.,
4 1988). River corridor dynamics are commonly studied through detailed observations at a
5 relatively limited spatial scale, which is restricted in an attempt to characterize local feedbacks
6 between mechanisms. These approaches often lack sufficient spatial resolution to enable
7 confident application of geostatistical approaches, and may not reliably support assessments of
8 system dynamics (e.g., Lee-Cullin et al., 2018). Longer-term dynamics are often considered as
9 historical contingencies: the biotic and abiotic histories or antecedent conditions that lead to the
10 present characteristics of the river corridor and affect its response to future perturbations.
11 Examples of river corridor studies that incorporate historical contingencies include perturbation-
12 response dynamics, commonly associated with floods (Czuba et al., 2019; Wu et al., 2018),
13 droughts (Boulton et al., 2004; Wood et al., 2010), or restoration activities (Rana et al., 2017;
14 Smidt et al., 2015), and large-scale historical perturbations such as land development (Liébault &
15 Piégay, 2002; Walling & Fang, 2003; Wohl, 2005), river regulation (Gregory, 2006), and
16 contamination (Byrne et al., 2012; Santschi et al., 2001). Such studies often involve little to no
17 replication and may be biased towards response variables that change rapidly relative to
18 processes that are quasi-steady over the timeframe of a given experiment.
19
20
21
22
23
24
25
26
27
28
29
30
31

32 While external forcing, internal dynamics, and historical contingencies have each been studied in
33 their own right, recent studies are beginning to integrate these concepts into holistic
34 understanding of river corridors. For example, Wisnoski and Lennon (2021) explicitly linked
35 localized heterogeneity to systematic spatial patterns along the network, revealing that the local
36 microbial assemblage in headwaters streams was controlled by local physical and chemical
37 conditions, but these local controls gave way to systemic organization from headwaters to larger
38 downstream rivers as the spatial scale of study increased. Such explicit consideration of local and
39 network scales is rare and still does not address historical contingencies. However, if done more
40 often and expanded to consider historical contingencies as a context for each replicate, this type
41 of systematic approach would allow assessment of the transition in dominant controls from local
42 heterogeneity (a reflection of internal dynamics) to larger-scale spatial organization (a reflection
43 of external drivers), the specific mechanisms of this transition, and the scale at which the
44 transition occurs, and how future multi-scale dynamics may depend on antecedent conditions (a
45 reflection of historical contingencies). Studies that have explicitly considered local
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 spatiotemporal dynamics as part of long-term system-wide functions have found strong
4 relationships between large-scale system structure, internal dynamics, and long-term emergent
5 outcomes in flow, sediment transport, and biogeochemistry (e.g., Fisher et al., 1998; Harvey &
6 Gooseff, 2015; Krause et al., 2017; Pinay et al., 2015). The success of these studies demonstrates
7 our ability to identify a core set of transferable and scalable processes that govern river system
8 dynamics and unify seemingly disparate observations into holistic understanding of river
9 corridor services and functions.

10
11 Here we use objective data-oriented approaches to confront the challenge of how a discipline
12 organized around the classic deductive model of site-, scale-, and mechanism-specific
13 observations can systematically link the resulting fragmented information into system-level
14 understanding. Our aim is to identify couplings that span scales and disciplinary expertise in
15 absence of pre-existing conceptual models that would traditionally serve as the source of
16 hypotheses for deductive testing. We propose an inductive approach to data synthesis, serving as
17 a basis for the unconstrained generation of new and potentially unexpected
18 hypothesesrelationships, each of which may be explained by hypotheses that could subsequently
19 be tested. To this end, we analyze a novel large data set for a 5th order river basin (Ward,
20 Zarnetske, et al., 2019) using inductive approaches to generate ~~novel~~a network of relationships
21 hypotheses that span traditional disciplinary boundaries. The data set contains 157 variables with
22 nearly 25,000 possible pairwise relationships, making it infeasible to explore each potential
23 causal pathwayrelationship through the lens of deductive inquiry. Further, the large degree of
24 covariation in environmental conditions may obscure underlying causal mechanisms, making it
25 difficult to determine unique process relationships and their controls. Thus, we pilot a machine
26 learning approach that sieves and categorizes information to identify non-obvious relationships
27 that merit subsequent investigation. We envision the apparent relationships generated by our
28 approach as a suite of observations around which hypotheses could be generated and
29 subsequently tested with more traditional approaches. In this way, we complement traditional
30 approaches by highlighting observations that may warrant hypotheses to be spun that explain
31 causal pathways that, thereby generating novel, interdisciplinary, and trans-scale to explain
32 hypotheses on river corridor dynamiesthe apparent relationships. This allows us to synthesize
33 complex, multi-scale observations independent of any pre-conceived conceptual models and

1
2
3 uncover novel and exciting information about the structure and function of river corridors. We
4 critically evaluate the resultant relationships relative to existing knowledge, and provide two
5 examples of how these novel insights may motivate future research questions that inform a
6 synthesis approach to understanding of river corridors.
7
8
9

10 11 12 **2. Methods**

13 **2.1 Data description and organization**

14 **2.1.1 Field site and synoptic campaign**

15
16 The H.J. Andrews Experimental forest (Western Cascades, Oregon, USA) is a 6,400 ha basin
17 that is primarily covered in old-growth and second growth forest and drained by a 5th order river.
18 The physical characteristics of the basin are well-described elsewhere (Deligne et al., 2017;
19 Dyrness, 1969; Jefferson et al., 2004; Swanson & James, 1975; Swanson & Jones, 2002). A
20 synoptic sampling campaign including detailed characterization of physical, chemical, and
21 biological characteristics and processes in the river corridor at 62 sites across stream orders 1-5
22 was conducted by Ward et al. (2019), which forms the basis of our study data set. These data are
23 the most uniform, comprehensive, and multi-scale available – to our knowledge – and, as such,
24 are optimal-~~uniquely useful~~ for ~~hypothesis-generation~~assessment of relationships spanning scales
25 and disciplines. Notably these data represent a spatial synoptic sampling design (i.e., a snapshot
26 in time), meaning their analysis will necessarily highlight apparent spatial patterns but cannot
27 capture the temporal dynamics of the system. Indeed, river corridors will have processes
28 operation spanning orders of magnitude in temporal scale (Ward and Packman, 2019).
29 Consequently, our approach will not capture temporal couplings between relationships, and we
30 are combining relatively dynamic variables (e.g., water temperature) and relatively static
31 variables (e.g., surficial geology) into a single analysis. -Approaches with comparable coverage
32 occurring through seasonal, storm, and/or diurnal fluctuations would enable a related assessment
33 of temporal dynamics and the persistence of relationships through natural variation.
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

50 **2.1.2 Data reduction**

51
52 Starting from this data set, we reduced the full suite of variables from Ward et al. (2019) to a
53 subset we considered to be most representative summary of the data set. For example, we
54 omitted identification of individual species and life-stages from macroinvertebrate data in favor
55
56
57
58
59
60

of summary indices, and similarly reduced ~~metabolomics data to a series of indices rather than attempting to explicitly analyze~~ 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

1
2
3 other than dependent variable as possible inputs, with the goal of identifying relationships
4 between individual variables. The second set used PC1 and PC2 from each disciplinary subgroup
5 as possible inputs with the goal of identifying more generalizable flows of information from the
6 major axes of variation within and between subgroups. In all cases SVMRs are used to identify
7 directional relationships between all possible pairs of variables (i.e., finding variable A is
8 informed by variable B does not require B is informed by A).
9

10
11
12
13
14
15 Finally, we compared performance of the SVMRs selecting features from the full variable set to
16 those selecting from a random subset. We constructed 100 SVMRs using 10 randomly selected
17 features as possible inputs for each variable. We used one-way ANOVA and Kruskal-Wallis
18 tests as a basis to assess performance differences between models with the full feature set vs.
19 random subset, reporting p_{ANOVA} and p_{KW} , respectively. We interpret SVMRs selecting from the
20 full feature set performing significantly better than those selecting from a random subset of
21 features as confirmation that the methods are identifying relationships that are at least
22 mathematically non-random.
23
24
25
26
27
28
29
30
31

32 **2.5 Literature analysis**

33
34 To assess the presence and relative frequency of studies jointly considering relationships
35 between each pair of variables in our data set, we conducted a series of searches using the
36 Scopus database in October 2020, following methods from similar studies (Ward, 2015; Yoder et
37 al., 2020). Each variable in our data set was assigned one or more keywords that are commonly
38 used to describe that variable in the literature (Ward, 2021). Literature was searched for every
39 pairwise combination of variables (12,246 unique searches) for studies containing both keywords
40 and a required term to indicate a study was likely relevant to our study of river corridors (one of:
41 river, stream, water, aquatic). We tabulated the total number of studies returned from each search
42 to assess the interactions between variables that have been studied jointly with greater or lower
43 frequency, and compared these results to the interactions found to be significant within the
44 SVMR analysis. Conversely, we also assessed if the specific pairwise interactions identified as
45 significant in the SVMRs were present in the literature.
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

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. 12A; 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

1
2
3 interpret PC1 as representing richness (high to low) and PC2 as representing abundance (high to
4 low). For stream solute tracers, we interpret PC1 as representing short-term storage of tracers
5 (low to high) and PC2 as representing the importance of advection and longitudinal dispersion to
6 tracer transport (low to high).
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

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

PCA on all variables						
	PC1			PC2		
	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
PCA on subgroups						
	PC1			PC2		
	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 5th 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

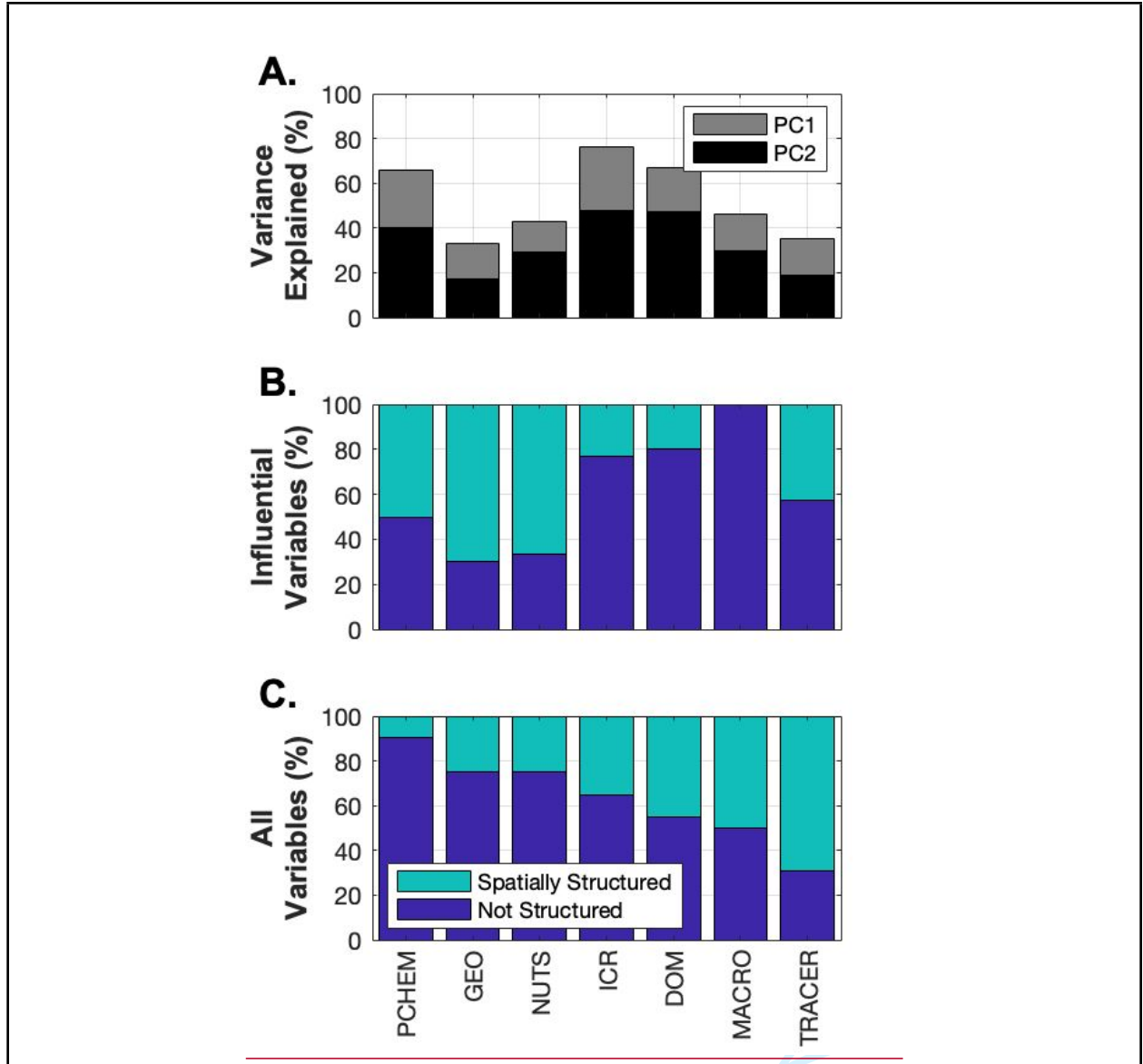
1
2
3 unconstrained in the longitudinal dimension until they reach the ocean. Still others have found
4 spatial structure in some parameters (e.g., in-stream solute concentrations) at scales that were
5 encapsulated within our study (e.g., McGuire et al., 2014), suggesting that finding of spatial
6 correlation lengths in one system or for one variable may not be universally transferable.
7
8
9

10
11
12 The fraction of influential variables with spatial structure was varied between subgroups (Fig.
13 1B, 1C), with 6 of 14 subgroup of PCs containing both structured and unstructured variables.
14

15 The largest proportion of spatially structured variables were in the nutrient-TRACER subgroup
16 (69%; Fig. 1C), and the least were in the macroinvertebrates-PCHEM subgroup (9.5%; Fig. 1C).
17

18 The variables that appear in the disciplinary subgroup PCs did not separate into distinct groups
19 of structured vs. unstructured variables. Instead, we found 44% of all influential variables were
20 spatially structured (23% in PC1 and 21% in PC2) compared to overall 36% of all variables
21 exhibiting spatial structure. All subgroups contained some structured influential variables except
22 for MACRO (Fig. 1B), where only unstructured variables were selected. Similarly, the fraction
23 of influential variables with spatial structure was consistent across subgroups (Fig. 1B, 1C), and
24 6 of 14 subgroup of PCs contained both structured and unstructured variables.
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



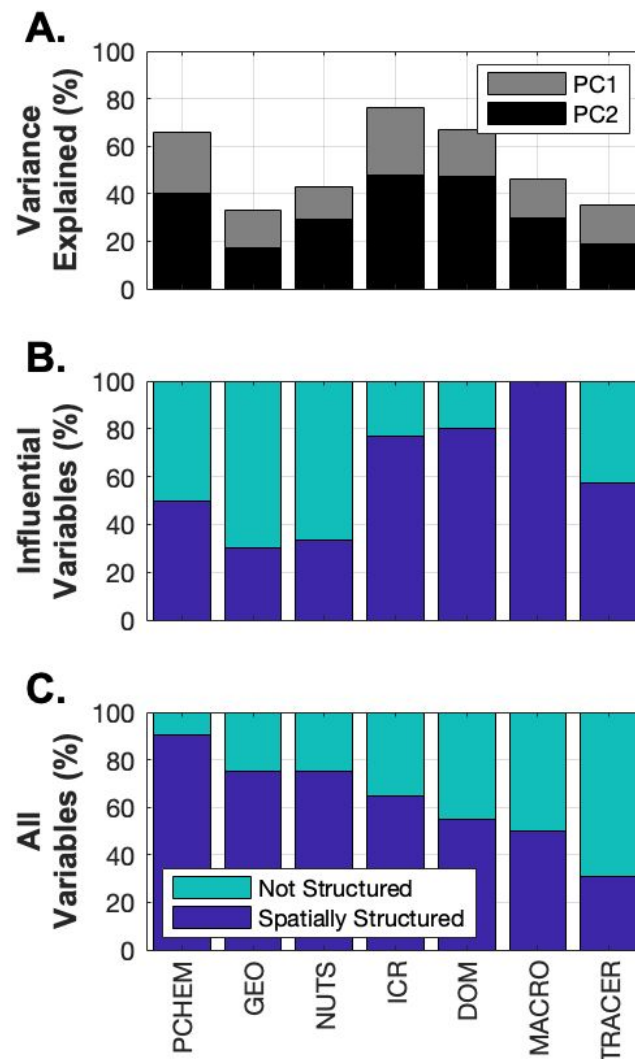


Fig. 1. (A) Variance in the Andrews river corridor data set explained by PC1 and PC2 for each expert subgroup. (B) Percentage of influential variables (i.e., the variables included in the first two PCs) that do and do not have spatial structure. (C) Percentage of all variables within each subgroup that 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 in the underlying data, with an overall mean r^2 of 0.83 (median 0.94, range 0.00 - 1.00). SVMRs

1
2
3 for individual variables selected an average of 4.4 variables as predictors (median 4, range 1 to
4 10; [Fig. S1](#)), indicating that the relationships (i.e., statistical models) identified by the SVMRs
5 were reasonably parsimonious. Additionally, performance of the SVMRs built from the full
6 feature set was significantly better than those built from a random selection of features ($p_{ANOVA} =$
7 $1E-19$; $p_{KW} = 4E-29$), indicating SVMRs are selecting meaningful features and the associated
8 relationships are appropriate for further analysis. The models built for spatially structured
9 variables had an overall mean r^2 of 0.91 (median 0.97, range 0.08 - 1.00) compared to a mean r^2
10 of 0.78 for unstructured variables (median 0.90, range 0.00 - 1.00). Goodness of fit was also
11 statistically better for the spatially structured variables ($p = 0.008$; one-way ANOVA), indicating
12 that spatially structured variables were more accurately predicted (i.e., higher r^2) compared to
13 unstructured variables.
14
15
16
17
18
19
20
21
22
23

24 Of the 157 variables predicted, 22% (34 variables) are informed by only out-of-group variables
25 (i.e., variables from a different subgroup), and 11% (17 variables) are informed by only within-
26 group variables (i.e., variables in the same subgroup). Thus, 67% of variables (106 out of 157)
27 required both in-group and out-of-group information for optimal prediction by the SVMRs.
28 Moreover, out-of-group information dominates predictor selection, representing an average of
29 59% of variables selected (median 66%, range 0-100%; Fig. 2, Table S1). Spatially structured
30 variables represent an average of 27.3% of variables selected for individual SVMRs (Fig. [S2A](#),
31 [S2C3](#)). Across the 157 SVMRs constructed, 30% (47 variables) did not select any spatially
32 structured features. We found 3% of models (5 variables) selected only spatially structured
33 features, and the remaining 67% (105 variables) selected a combination of structured and
34 unstructured variables.
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

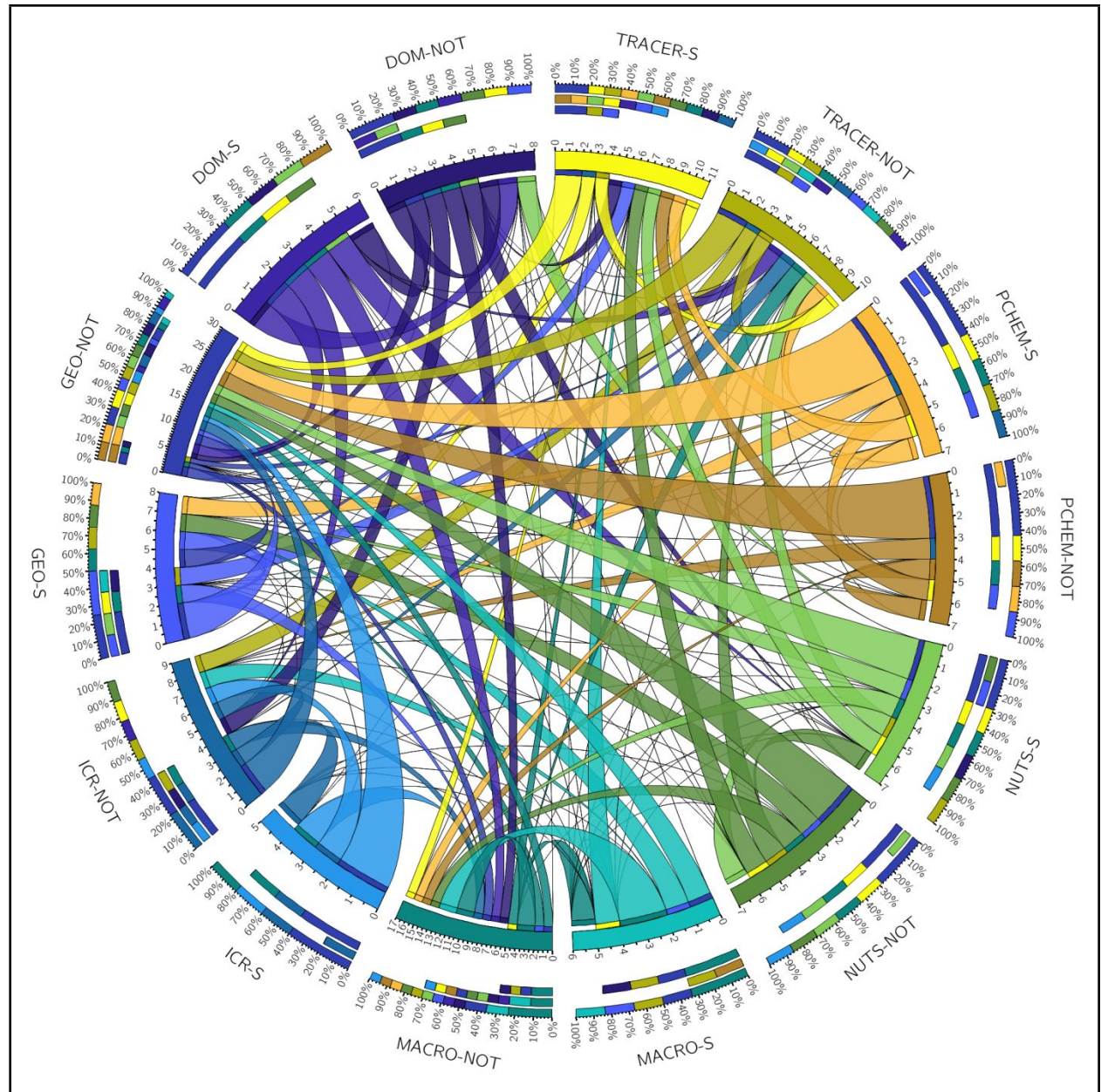


Fig. 2. Information flow within and among subgroups of variables commonly used as measures of river corridor dynamics based on the suite of SVMs constructed for each variable (Section 3.3.1). The variables included in the 7 subgroups (PCHEM = physical chemistry; GEO = geologic setting; NUTS = nutrients; 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 variables subgroups as:

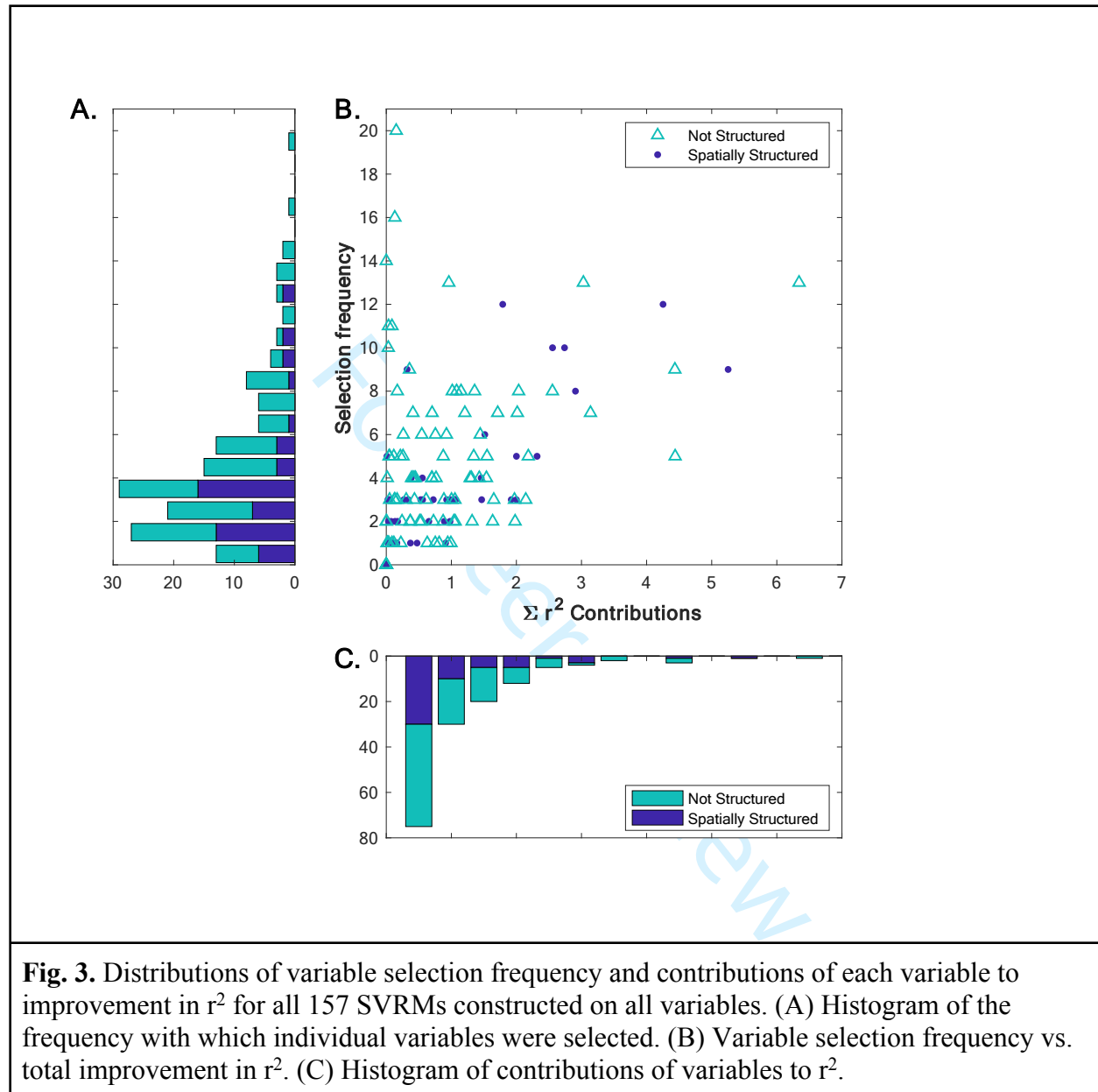
- Inner Ring: destination(s) of information from each subgroup (i.e., answers the question "which other subgroups used information from this subgroup?"; colloquially the 'outflows' of information from one subgroup to another). These are the independent variables requires as inputs to make predictions of dependent variables in other groups.
- Inner-Middle Ring: the source(s) of information to a subgroup (i.e., answers the question "which variable groups 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. contributed information to the predictions for the given group).
- Middle Ring: destination of information from each subgroup (i.e., which groups needed information from a given group for their predictions). Outer Ring: Scaled, tTotal 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?" the sum of the inner and middle rings). 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 NH₃ concentration. However, this variable only contributed 0.046 improvement in r² 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 r² summed across the 20 models where it was selected (average improvement of 0.315 in r² 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 r² of the model. We found out-of-group variables represent about 30% of all selections within the SVMRs ([Fig. S2Ce](#)), but contribute more than 50% of the improvements in model performance ([Fig. S2Dd](#)). 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

1
2
3 more frequently than PCs from any other subgroup (Table S2). In fact, about 25% of variables
4 (39 of 157) were predicted solely from their in-group PCs. The explanatory power of PCs for in-
5 group variance is unsurprising given that PC1 and PC2 were successful in explaining an average
6 of 52% of variance within their group. However, we also found about 26% of variable
7 predictions (41 of 157) used only out-of-group PCs, and 118 variable predictions selected at least
8 one out-of-group PCs. Further, variables in each subgroup drew information from nearly every
9 other subgroup (see Table S1), These findings indicate that studies that are limited to one
10 discipline are unlikely to explain as much of the observed variance in the measured variables as
11 studies that intentionally span disciplinary boundaries, and that it is important for disciplinary
12 understanding to at least characterize the major attributes from other subgroups.
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

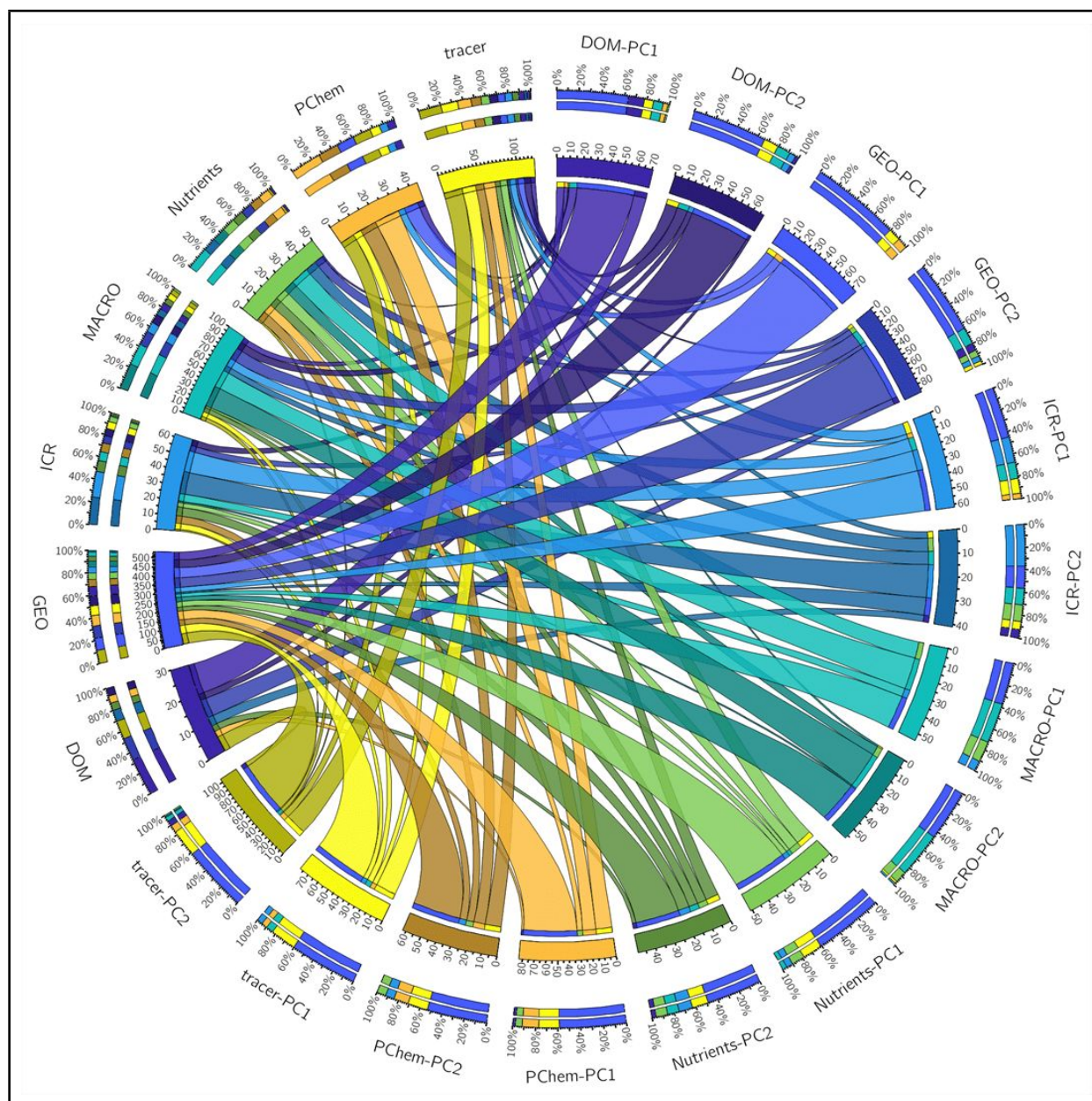


Fig. 4. Circos plot showing the one-way flow of information from the subgroup PCs (Table 1; labeled “XXX-PCY” where XXX is the subgroup and Y in the PC number) to variables predicted by the suite of SVMRs described in Section 3.3.2. Plot layout and interpretation is identical to that described for Fig. 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

1
2
3 Our literature search identified 4,075 combinations of variables that have been studied pairwise
4 in the literature (of 12,246 possible combinations). The pairwise literature search returned a total
5 of 2,731,694 results. The number of studies identified for any given pair of variables was highly
6 skewed: 50% of published studies included the 18 most commonly studied pairs of variables
7 [\(Ward, 2021\)](#), while the number of studies of any given pair of variables ranged from 1 to
8 270,015 (mean 670, median 14). These findings indicate a bias toward co-observation and
9 reporting of a limited number of pairwise studies, consistent with a prior study that manually
10 reviewed search results (Ward, 2015). We also found the existing literature is more focused on
11 in-group relationships (57.2% of pairwise results) compared to between-group relationships
12 (42.8% of pairwise results). In contrast, our SVMR approach identified a total of 672 pairwise
13 relationships, of which 68.8% are between-group. Notably, about 84% or 564 variable pairs do
14 not appear to have been reported previously (i.e., our systematic literature search did not return
15 any manuscripts containing information on both variables). The remaining ~~28.216~~28.216% (108
16 relationships) have been previously reported in the literature (Fig. 5; [Fig. S5](#); [Table S4](#)). The 108
17 relationships found in both the literature and in our data analysis only represent about 2.6% of all
18 previously-reported relationships, but these relationships are included in more than 16% of all
19 published studies, indicating that prior studies have focused primarily on a relatively small
20 number of relationships. On the basis of within- and between-group frequency, the literature is
21 broadly not reflective of our findings, with the SVMR identifying higher frequencies of between-
22 group relationships that are present in the literature (Table S3). Finally, we note that the lack of a
23 relationship in the SVMR does not necessarily indicate that some relationship may be possible,
24 just as the presence of a statistical relationship does not necessarily indicate a causal relationship.
25 Some meaningful relationships could have been omitted due to signal-to-noise ratios, lagged
26 correlation between variables, or because a highly correlated variable was already selected. This
27 may explain why some well-studied relationships were not apparent in our analysis (Fig. 5).
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

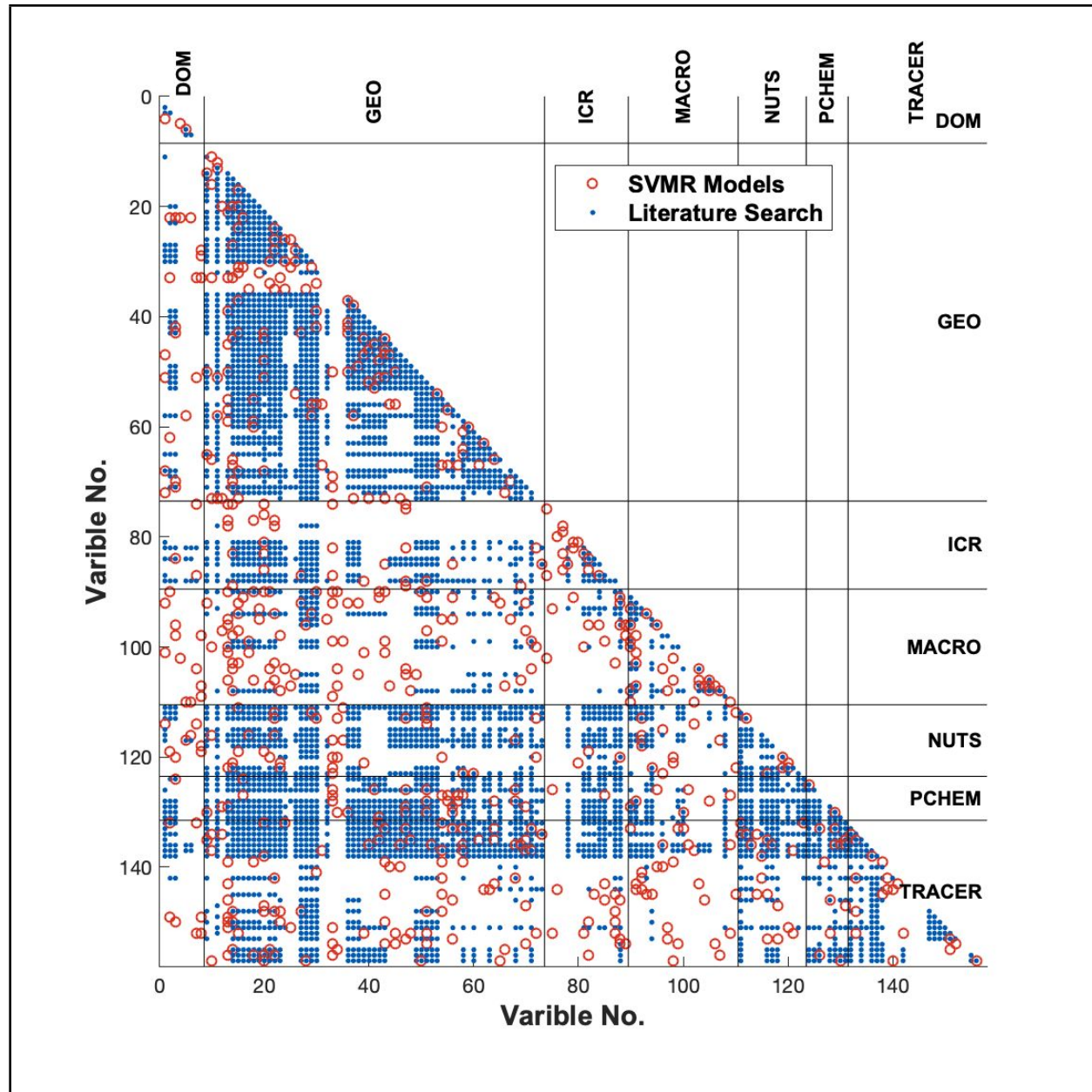


Fig. 5. Scatterplot showing pairwise study in the literature (blue dots) and identification of a relationship in our SVMR approach (red circles) for all variable pairs. Variable numbers correspond to the order variables are listed in Table S1.

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

1
2
3 27% of the variables included by SVMRs although they make up 36% of the total variable set).
4 This means the predictions of spatially structured variables were not dominated by structure from
5 a small number of structured variables. Further, a majority of variables observed (about 64%)
6 were not themselves spatially structured, and five of the seven subgroups (PCHEM, GEO,
7 NUTS, ICR, TRACER) resulted in at least one PC that was not spatially structured. These results
8 indicate that spatial structure is not ubiquitous in the river corridor. Instead, some variables
9 represent local ‘noise’ on the network-scale ‘signal’ (i.e., systematic variation in physical,
10 chemical, and biological processes from headwaters to large rivers; Vannote et al. 1980). This
11 heterogeneity is either independent from large-scale system structure (i.e., controlled by local
12 process interactions that are neither controlled by nor influence larger scale patterns) or simply
13 have sufficiently high variability to obscure larger-scale trends. Such localized ‘noise’ may also
14 reflect processes whose importance is localized in space or time, but do not recognizably follow
15 a larger spatial structure.
16
17
18
19
20
21
22
23
24
25

26
27 Individual variables reflect complex interactions that can either lead to the emergence of spatial
28 structure or overwhelm the underlying spatial structure associated with more basic variables like
29 slope and elevation. We found six variables that were spatially structured but had strong
30 relationships (SVMRs) that only included unstructured variables. In these cases, spatial structure
31 emerged or was generated by the interaction of variables that did not themselves have spatial
32 structure. Conversely, 60 of the SVMRs for unstructured variables included at least one spatially
33 structured variable (38 selected 1, 14 selected 2, and 8 selected 3 spatially structured variables).
34 This pattern suggests that spatial structure does not necessarily propagate from one variable to
35 another, indicating “signal shredding” in the river corridor (Jerolmack & Paola, 2010), where
36 information is erased by interactions between variables. While such behavior has only been
37 confirmed previously for sediment transport, our findings indicate that localized feedbacks can
38 generally overwhelm underlying spatial structure within the river corridor. This suggests that
39 sufficiently large perturbations will have system-wide impacts (e.g., large fires, floods), but
40 internal dynamics may overwhelm large-scale patterns under normal circumstances.
41 Consequently, studies of river corridors must consider local-scale interactions (i.e., internal
42 dynamics), large-scale drivers (i.e., external forcing), and the temporal context (i.e., historical
43 contingencies) if we are to account for the feedbacks and interactions in the river corridor.
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

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 as well as 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 can 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

1
2
3 observations (Payn et al., 2009; Valett et al., 1996) and models (Cardenas, 2008; Frissell et al.,
4 1986; Wondzell & Gooseff, 2014; Wörman et al., 2007).
5
6
7

8 Ward et al.'s (2019) observation of monotonic trends between most hydrologic exchange metrics
9 and discharge - which they describe as a proxy for network position - agree with our finding of
10 spatial structure in several variables describing geomorphic setting (including hydraulic
11 conductivity, valley slope, valley width, sinuosity), river flow (velocity, discharge), and solute
12 transport metrics (e.g., median travel time, skewness). We did not find spatial structure for other
13 metrics of exchange where Ward et al. did, including the coefficient of variation, holdback, and
14 channel water balance. Further, many of the relationships identified by Ward et al. have low
15 explanatory power as evidenced by low r^2 values, indicating that hydrologic exchange cannot be
16 described by a single explanatory variable. In contrast, the multivariate and nonlinear responses
17 encoded in the SVMs better explain the patterns in river corridor exchange observed in the
18 Andrews watersheds.
19
20
21
22
23
24
25
26
27
28

29 **4.2.2 Microbial Community Assembly**

30 Interactions along the river corridor can not only 'shred' or erase information (*sensu* Jerolmack
31 & Paola, 2010), but can also generate new information and patterns. For example, prior work at
32 the H.J. Andrews Experimental Forest spanning headwaters through 5th order rivers (Wisnoski
33 and Lennon, 2021) showed that microbial assemblages in headwater streams ~~were~~-habitat-
34 dependent, while the microbial community became more homogeneous with distance
35 downstream. Additionally, ~~Wisnoski and Lennon found that~~the same study found taxonomic β -
36 diversity was explained by an axis with positive loadings for elevation and dissolved organic
37 carbon, and negative loadings for electrical conductivity, pH, total nitrogen, and total
38 phosphorus. Microbial assemblages are known to arise in response to local heterogeneity in the
39 landscape, integrating inputs and environmental variables in space and time. While we did not
40 analyze microbial assemblages explicitly here, we can interpret our observations in the context of
41 prior findings at the site (Wisnoski and Lennon, 2021). ~~do compare geomorphic and water~~
42 ~~quality variables with prior observations of the microbial community assemblage.~~ Our results
43 show spatial structure in electrical conductivity and several geomorphic variables that are known
44 to vary with elevation, but no spatial structure in total dissolved phosphorus, DOC, or total
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 dissolved nitrogen. In comparison to the controls on taxonomic β -diversity described by
4 Wisnoski and Lennon (2021), we did find spatial structure in elevation, in-stream nitrate+nitrate,
5 and electrical conductivity, but not in bulk dissolved organic carbon, ammonia, or total
6 phosphorous. Thus, our findings are broadly consistent with past findings that at least some of
7 the known controls on microbial diversity are spatially structured. However, we also note that
8 not all controls were structured, but the related microbial community did retain spatial
9 organization. Thus, we interpret the spatial organization of the microbial assemblage as the
10 emergence of spatial structure from a suite of largely unstructured variables in the river corridor.
11
12 Consequently, studies focused at single locations along a stream may be missing contextual
13 information on controlling factors that have propagated from the catchment headwaters, or
14 misinterpreting signals that were generated within the river corridor itself.

24 **4.2.3 River Continuum Concept**

25 The River Continuum Concept (Vannote et al., 1980) -- one of the most widely recognized and
26 cited conceptual model of river corridors -- argues that Leopold's conceptual model that
27 geomorphology reflects energy equilibrium can be extended into ecosystem functions (Langbein
28 & Leopold, 1966; L B Leopold et al., 1964; Luna B. Leopold & Langbein, 1962). Vannote et al.
29 (1980) specifically proposed: (a) biological communities should evolve to optimize the use of
30 available energy (i.e., biodegradable organic matter); and (b) energy availability will vary
31 systematically from headwaters to large downstream rivers. Our PCA results on all variables are
32 broadly consistent with these hypotheses, which is to be expected at the H.J. Andrews
33 Experimental Forest was one of the key sites studied in developing and demonstrating the
34 conceptual model. We found organic matter -chemistry chemistry and geological setting
35 explained 37% of the variance across the entire data set (PC1 and PC2; Table 1). With regard to
36 biological communities optimizing to use available energy in an organized fashion, we do see
37 that available energy – in this case assessed via energy availability in organic carbon (PC1 on all
38 variables) – defines one critical dimension of variation in the system. Additionally, the high
39 proportion of spatially organized variables in TRACER, GEO, and NUTS is consistent with
40 broad concepts of systematic organization along river networks. Indeed, we ~~We also~~ found
41 spatial structure in about 36% of all variables across all disciplinary subgroups, consistent with
42 the idea that large-scale gradients will drive systematic trends in both physical and

1
2
3 biogeochemical processes. We did find spatial organization in shredders which is consistent with
4 the River Continuum Concept. ~~Six of the fourteen subgroup PCs were spatially structured~~
5 ~~(Table 1), reflecting broad spatial structure in the H.J. Andrews catchment. Our findings of broad~~
6 ~~patterns along the river network, as evidenced by spatial structure, is broadly consistent with the~~
7 ~~River Continuum Concept, which was based on a much more limited set of measurements.~~ Our
8 findings on the importance of organic carbon as an explanatory variable for patterns in the river
9 corridor also support Vannote et al.'s expectation of the importance of energy availability to the
10 structure of fluvial ecosystems.
11
12
13
14
15
16
17

18 **4.3 ~~Novel hypotheses and~~ Open questions stemming from the inductive analysis**

19 We applied machine learning techniques to cross-disciplinary data to uncover novel **hypotheses**
20 **relationships** that are worthy of subsequent investigation. Inductive approaches cannot reveal
21 causal relationships, making this a useful approach to identify relationships for future study,
22 rather than proving mechanistic pathways. To demonstrate the value of this approach, we explore
23 a selection of findings from the network of relationships identified by our SVMR models,
24 focusing on relationships that ~~have not been previously identified and are not likely to be~~
25 ~~uncovered or explored through conventional approaches~~ are at the cutting edge of our
26 understanding of river corridors. While our body of knowledge has methodically built
27 knowledge and is beginning to engage with these questions, we take it as a positive sign that
28 inductive approaches were able to also pick these relationships out of the data set. Thus, in
29 addition to consistency we past findings (Section 4.2) we take these findings as further support
30 that inductive approaches are able to identify relationships worthy of further scrutiny. We pose
31 these as hypotheses-potential areas for future study to highlight the role of inductive analysis as a
32 path to inspire the asking of questions, rather than providing mechanistic -answers, questions
33 about the complex structure and function of river corridors.
34
35
36
37
38
39
40
41
42
43
44
45
46
47

48 **4.3.1 Why are metabolomics data most informed by geological variation?**

49 Metabolomics data alone formed PC1 for the overall analysis, explaining 20% of the variation in
50 all data analyzed (Table 1), while geomorphic variables dominate PC2, explaining 17% of all
51 variance. Moreover, these axes are, by definition, orthogonal implying that the two groupings
52 should be independent. Across the 16 SVMRs constructed on organic carbon chemistry (ICR)
53
54
55
56
57
58
59
60

1
2
3 variables, none selected any features from the dissolved organic matter, nutrient, nor physical
4 chemistry subgroups (DOM, NUTS, and PCHEM, respectively). Instead, out-of-group
5 information was exclusively from geological features, solute tracer, and macroinvertebrate
6 groupings (GEO, TRACER, and MACRO, respectively). This is particularly surprising given
7 that a host of variables traditionally used to describe organic matter were available, including
8 optical measures of carbon quality (e.g., EEM features, SUVA₂₅₄) and quantity (e.g., total DOC,
9 carbon acquiring extracellular enzymes). We posit that the apparent dominance of physical
10 setting over biogeochemical variables emerges through the microbial community (i.e., the Baas
11 Becking hypothesis; *sensu* O'Malley, 2008; Fondi et al., 2016; Wit and Bouvier, 2006). In other
12 words, geologic setting and hydraulics set a template that defines which microbial communities
13 will occur, and these communities are responsible for the molecular form of organic matter that
14 is transformed within and exported from a given location. This is, functionally, the River
15 Continuum Concept applied to microbial communities. We expect the role of microbial
16 community structure in defining ecosystem processes will be critical as we transition from
17 conceptual models based on bulk measurement of organic matter (e.g., DOC, EEMs) to models
18 informed by metabolomics.
19
20
21
22
23
24
25
26
27
28
29
30

31
32 Previously developed theories based on bulk DOC or proxies for organic matter quality must be
33 revisited, because the field of metabolomics is rapidly evolving. The limited suite of studies that
34 include both organic carbon chemistry and nutrient data (ICR and NUTS) make comparisons for
35 consistency of findings limited. It is possible that previous conclusions about carbon limitations
36 in some systems may have been biased by only considering bulk DOC or DIC instead of its
37 molecular composition, which is highly nonuniform in its ecological function. We do not expect
38 that organic matter molecular composition is entirely controlled by geologic setting (though such
39 control has been reported; e.g., Robertson et al., 2019; Cotrufo et al., 2013), but instead that in-
40 stream organic matter reflects the integration of physical, chemical, and biological processes
41 occurring upstream of the sampling location. These processes are diverse, spanning the
42 influences of terrestrial vegetation, soil-forming processes, photochemistry, organo-mineral
43 interactions, and in-stream biological production and transformation of organic molecules. Thus,
44 the core questions are to understand when, where, and how organic matter is produced,
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 transformed, and transported. We expect that understanding microbial communities and their
4 metabolism will be critical to answering these questions.
5
6
7

8 In addition, Danczak et al. (2020) proposed a conceptual framework that draws parallels between
9 organismal birth, death, and dispersal and organic matter production, transformation, and
10 transport. They argue that organic molecules are assembled into metabolomes via a combination
11 of production, transformation, and transport just as organisms are assembled into communities
12 via a combination of birth, death, and dispersal. Danczak et al. (2020) also provide an analytical
13 approach for quantifying assembly processes, including the ability to infer when transport
14 overwhelms influences of production and transformation. This approach may be fruitful in
15 linking upland dynamics to aquatic dynamics (Waring et al., 2020; Wisnoski et al., 2021),
16 linking microbial community assembly processes to organic matter assembly processes, and
17 further highlights the need for conceptual synthesis in the river corridor (Stegen et al., 2018).
18
19
20
21
22
23
24
25
26

27 ~~Further~~Finally, metabolomics data has been used previously to inductively reveal limitations of
28 using bulk water chemistry in river corridors to understand specific biogeochemical conditions.
29 For example, there has been a recent revelation that conceptual models for denitrification in river
30 corridors were framed at a large river network scale and not capturing dynamic, small scale
31 controls of anaerobic metabolic pathways, including denitrification (e.g., Briggs et al., 2015).
32 Since this revelation, field experiments and deductive methods have revealed that denitrification
33 is in fact occurring in sediment “microzones” across a wide range of river corridor conditions
34 that was previously hidden by and assumed impossible based upon bulk water chemistry (e.g.,
35 Knapp et al., 2017; Hampton et al., 2019; Hampton et al., 2020).
36
37
38
39
40
41
42
43
44

45 ***4.3.2 What controls nitrogen-acquiring extracellular enzymatic activity in a nitrogen-limited*** 46 ***ecosystem?*** 47

48 Aquatic ecosystems at the H.J. Andrews have been historically considered to be nitrogen limited
49 (Sollins et al., 1981; Triska et al., 1984). Consequently, we expected that microbes would
50 generate both leucine aminopeptidase (LAP) and N- acetylglucosaminidase (NAG) to acquire
51 nitrogen and that this would be ubiquitous across the basin. Moreover, C:N:P ratios of
52 extracellular enzymatic activity (EEA) should indicate an overproduction of N-acquiring
53
54
55
56
57
58
59
60

1
2
3 enzymes as N-limited microbes allocate energy to acquiring their limiting nutrient (e.g.,
4 Sinsabaugh et al., 1997) .
5
6
7

8 To test this expectation, we considered two nitrogen-acquiring enzymes: LAP and NAG. LAP
9 was part of PC1 for the NUTS subgroup and was orthogonal to total organic matter in the
10 sediment, indicating little control on sediment organic matter in explaining LAP. SVMRs for
11 LAP identify several GEO variables (bedrock type, hillslope stability, and channel water
12 balance), allochthonous inputs to the river (deciduous forest, abundance of collector-gatherer
13 macroinvertebrates), and organic carbon (spectral slope and ICR 'other molecules'). Positive
14 correlations with spectral slope and small molecules in the ICR indicate increased LAP occurs
15 where relatively small and non-aromatic carbon sources are present. Similarly, NAG was
16 predicted by bedrock type, ICR (protein abundance), and phosphorus-acquiring enzymes.
17 Because we do not see spatial structure in LAP, NAG, nor 11 of the 13 variables selected by
18 their SVMRs, we infer that there is not a spatial control on nitrogen acquiring enzymes.
19
20
21
22
23
24
25
26
27
28

29 Several studies have reported increasing EEA with nutrient availability (Hill et al., 2010;
30 Sinsabaugh et al. 1997; Williams et al. 2010; Williams et al. 2012), which is not consistent with
31 our findings (i.e., no measurement of bulk nitrogen, carbon, phosphorus, nor oxygen were
32 selected by SVMRs for the ICR subgroup). Instead, we find that EEA may be explained by
33 particular classes of organic matter – specifically smaller, less aromatic carbon molecules,
34 consistent with Williams et al. (2012) and Hill et al. (2010). We also hypothesize the prevalence
35 of GEO features selected by SVMRs but lack of spatial structure may indicate that there are
36 geogenic micronutrient controls on the localized enzymatic activity that have not been measured,
37 such as the availability of potassium, manganese, iron, and silica that weathers from local features.
38
39
40
41
42
43
44
45

46 Another enzymatic question that requires more deductive work is whether the entire river
47 corridor is N-limited. Ecoenzymatic ratios of 1:1:1 C:N:P suggest an equilibrium between
48 microbial biomass and detrital organic matter (Sinsabaugh et al., 2009). The ratios of C:N and
49 C:P acquiring enzymes in our study (GLU:LAP+NAG and GLU:AP, respectively, based on data
50 in Ward et al., 2019) have slopes that are statistically indistinguishable from analyses of global
51 datasets (Sinsabaugh and Shah, 2012), indicating EEA is produced in relative proportions to the
52
53
54
55
56
57
58
59
60

1
2
3 basic C:N:P ratios required by microbes, suggesting that the sediment microbial community may
4 not, in fact, be N-limited relative to the availability of other nutrients and substrates. Therefore,
5 while catchment-scale mass balances indicated one understanding of the system as N-limited
6 (e.g., Sollins et al., 1981; Triska et al., 1984), we interpret the EEA data as an indicator that the
7 microbial community has adapted to the available N, and that this is present across the network
8 (based on the lack of spatial structure).
9
10
11
12
13
14

15 Our analyses suggest many fruitful paths forward for interdisciplinary river corridor research.
16 These include, but are not limited to, the examples presented above that (a) relate molecular
17 characterization of carbon to EEA to investigate organic matter quality controls; (b)
18 comprehensively sample stream, streambed sediment, hyporheic pore water, and hyporheic
19 sediment communities for EEA to test our hypotheses that microbes are not N limited across
20 these spatial domains; and (c) use repeated measurements to assess if one spatial snapshot of the
21 network adequately captures temporally dynamic behavior (as was found in Giraldo et al., 2014).
22 Our findings also suggest that the concept of ecological stoichiometry and nutrient limitations
23 manifest differently across multiple scales, warranting consideration of the places, times, and
24 scales at which equilibrium or limitation should be inferred, and whether findings of limitations
25 at one scale can be directly transferred to other scales. One particularly compelling question
26 resulting from our work is whether system-wide, large-scale N-limitation indicate low N inputs
27 at all scales, internal limitations due to spatial structure or heterogeneity (e.g., localized inputs
28 from N-fixing alders), biogeochemical limitations (e.g., kinetics of organic matter breakdown),
29 or transport limitation (e.g., inaccessibility of nutrients in some locations)?
30
31
32
33
34
35
36
37
38
39
40
41
42

43 **4.4 Inductive relationships are hypotheses-observations around which hypotheses can be** 44 **spun and tested that warrant additional scrutiny**

45
46 The suite of models we constructed include 672 apparent relationships, 84% of which have not
47 been previously studied based on our literature search. It is important to recognize the
48 relationships identified here are intended as future directions, not as endpoints that reflect a
49 causal or mechanistic understanding, particularly in the case of correlations that have not been
50 reported by other studies. Each relationship serves as a set of observations, the first step in the
51 scientific method. We envision the next step for each relationship being the generation of
52
53
54
55
56
57
58
59
60

~~hypotheses that propose mechanisms or explanations, followed by rigorous investigation with deductive approaches~~ ~~t must be considered in the context of hypothesized mechanisms or explanations, and rigorously tested~~ 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 ~~each~~ hypotheses that can explain each inductive relationship as a hypothesis, 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 hypotheses-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 highly-focused 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

1
2
3 apparent in the initial study design. In this context, the inductive approach we propose here is
4 extremely useful for rapidly identifying relationships spanning disciplinary boundaries that
5 would otherwise take decades of disciplinary inquiry to identify.
6
7
8
9

10 ***4.4.2 Existing data sets are incomplete and could not have uncovered relationships***

11 Our analysis relies on the most comprehensive catchment-scale observations of interacting
12 physical, chemical, and biological processes in any river corridor to-date. The dataset we
13 analyzed also builds upon extensive prior work and data from the H.J. Andrews Experimental
14 Forest. Such comprehensive datasets, particularly co-located with long term ecological research,
15 have not previously been available and require extensive interdisciplinary collaboration to
16 obtain. For example, molecular organic matter chemistry (e.g., FTIRCMS) is only recently
17 emerging as part of river corridor science (Graham et al., 2018; Stegen, Johnson, et al., 2018;
18 Zhou et al., 2019) and has not been jointly collected with the breadth of observations we
19 analyzed here. To make further progress in unraveling the complexity of river corridors, we
20 recommend combining standardized system characterization across many streams and rivers with
21 intensive study of select watersheds to generate the rich datasets needed to evaluate process
22 interconnections and scale dependencies (Stegen & Goldman, 2018). In this case, the
23 comprehensive nature of the data set explains why novel relationships were identified here: such
24 breadth of data were simply not collected in past efforts. This further demonstrates the utility of
25 inductive analysis in generating hypotheses from new datasets that can then be tested more
26 broadly. Finally, note that our own data set, while comprehensive, is far from complete in terms
27 of all variables that could be measured across all relevant spatial scales, temporal scales, and
28 process dynamics.
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44

45 ***4.4.3 Relationships may be scale- or time-dependent***

46 Both the structure and function of river corridors are known to be scale-dependent (Frissell et al.,
47 1986; Rodríguez-Iturbe & Rinaldo, 1997; McCluney et al., 2014). The network scale considered
48 here is larger than many studies of river corridors (see reviews by Tank et al., 2008; Ward,
49 2015). It is possible that the relationships identified between variables here by SVMR do not
50 hold at all scales, or that the relationships are real but have not been tested over the range of
51 scales we included in our analysis. Prior studies of river structure have found that self-
52
53
54
55
56
57
58
59
60

1
2
3 similarities and scale dependencies generally only occur over a limited range of scales, and either
4 average out at large scales or are limited by a physical constraint (e.g., water depth, channel
5 width, valley width) (Jerolmack & Paola, 2010; Nikora & Hicks, 1997; Rodríguez-Iturbe &
6 Rinaldo, 1997). As with relationships between individual variables, scale dependencies and
7 scaling limits identified from broad data analysis must be considered as hypotheses and tested
8 using directed observations and/or simulations with competing or alternative formulations.
9 Similarly, analyses here focused on a data set collected under baseflow conditions and process
10 controls are expected to vary in response to seasonal and storm dynamics in forcing. Moreover,
11 our analysis are focused on what can be gleaned from a single snapshot in time, whereas the
12 actual characterization includes a combination of variables spanning relatively dynamic (e.g.,
13 dissolved oxygen) to relatively static (e.g., valley slope), which may cause some relationships to
14 manifest and obscure others. Future efforts to combine high temporal resolution data with spatial
15 synoptic campaigns could directly address this limitation.
16
17
18
19
20
21
22
23
24
25
26

27 ***4.4.4 Spurious correlation may have driven the inductive relationships identified***

28 The relationships identified in our study may represent spurious correlation of disparate data or
29 other mutual dependencies in the underlying data, a known limitation of machine learning
30 approaches. In this case, the inductive approach aids in identifying mathematical artifacts rather
31 than causal pathways or process interactions. Such relationships could also reflect redundant
32 information (i.e., several different variables may reflect similar features on the landscape, and the
33 autocorrelation amongst independently-measured variables may obscure underlying
34 relationships). For example, if geology, land cover, and soils all systematically vary with
35 increasing elevation, then these variables will all show consistent relationships that may
36 confound interpretation. We emphasize here the relationships identified by SVMR and other
37 machine learning methods only provide a starting point for generation of hypotheses, not an
38 endpoint. The next step for investigation of such putative relationships would be to hypothesize a
39 causal mechanism and design a study to collect the specific data needed to test it, while still
40 capturing the essential system information identified here for purposes of evaluating scale
41 dependency and complex system controls.
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

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

1
2
3 into spatially and temporally comprehensive datasets (e.g., Tiegs et al., 2019; Stegen and
4 Goldman, 2018).
5
6
7

8 In this study, we have demonstrated an application of machine learning approaches to generate
9 ~~hypotheses-relationships~~ that may ~~inspire new studies to ultimately reveal~~ ~~serve as~~ the
10 ‘connective tissue’ ~~that linking~~ our understanding across spatiotemporal scales and disciplines.
11 Indeed, the step of organizing raw observations to develop testable hypotheses is at the core of
12 the scientific method, ~~and we have prototyped one approach to organize observations and~~
13 ~~highlight potential relationships in the data~~. Hypothesis generation is touted as one of the core
14 values of field-based observation and monitoring (Burt & McDonnell, 2015; Lovett et al., 2007),
15 where observations demand explanations. The inductive approach used here presents a body of
16 putative relationships for subsequent study, at least some of which are consistent with prior
17 conceptualizations and observations of river corridors (~~i.e.,~~ section 4.2) ~~and emerging areas of~~
18 ~~inquiry (section 4.3)~~. We do not propose that such approaches supplant deductive science, but
19 rather that the two approaches must be coupled in river corridor science. The inductive approach
20 provides an unbiased or naive data synthesis, which has the potential to reveal patterns and
21 relationships that would not be obvious from our present, disciplinary perspectives.
22
23
24
25
26
27
28
29
30
31
32
33

34 5. Conclusions

35 We began with the assumption that all variables may interact with all other variables, yielding
36 nearly 25,000 ~~hypothesized-potential pairwise~~ relationships ~~between variables~~. Using machine
37 learning, we rejected most of these ~~hypotheses-relationships~~, identifying 672 ~~pairwise-apparent~~
38 relationships that ~~have explanatory power in the data set could not be rejected by this approach~~,
39 notably including 564 pairwise relationships that were not previously explored in the literature.
40 Put another way, we have generated a web of 564 new ~~hypotheses-apparent relationships~~ that
41 may reveal new couplings in the river corridor. These relationships eschew disciplinary or
42 method-specific approaches, providing ‘connective tissue’ between traditional discipline-, scale-,
43 site-, or method-dependent knowledge. Moreover, the network of relationships we have
44 identified is consistent with several past studies from the field site (Vannote et al., 1980; Ward,
45 Wondzell, et al., 2019; Wisnoski & Lennon, 2021), providing confidence that at least some of
46 these relationships are more than spurious correlations.
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5 Most of the relationships we identified, including a majority of those not present in the literature,
6 include between-group flows of information. Our results show that interactions between
7 processes that are typically studied by different disciplines is critically important to explain
8 structure and function in the river corridor. This conclusion is, perhaps, unsurprising as a
9 macrosystems view would acknowledge and expect to find cross-scale and interdisciplinary
10 relationships (Heffernan et al., 2014; McCluney et al., 2014). Still, this view is seldom fully
11 captured in existing experimental designs and the resulting data sets and literature. Importantly,
12 we also demonstrated that spatial structure can be both generated through the interaction of
13 unstructured data as well as destroyed or overprinted along the network. Thus, consideration of
14 how an observed pattern may emerge or not be visible along a spatial gradient is a critically
15 important consideration prior to interpretation of data sets.
16
17
18
19
20
21
22
23
24

25 Building connections between existing studies requires explicitly planning for synthesis in future
26 efforts. Here, we demonstrated the value of collecting data sets that enabled synthesis within and
27 between locations, disciplines, and scales. This does not diminish the value of traditional,
28 disciplinary hypothesis testing and deductive approaches to science. Instead, common metadata
29 and even a small number of out-of-group observations may enable synthesis efforts based on
30 inductive approaches that aids in spinning new hypotheses. Ultimately, inductive approaches are
31 a useful way to generate hypotheses from existing observational datasets and advance our
32 scientific understanding.
33
34
35
36
37
38
39
40

41 **Acknowledgements.**

42 This research has been supported by the Leverhulme Trust (Where rivers, groundwater and
43 disciplines meet: a hyporheic research network), the UK Natural Environment Research Council
44 (grant no. NE/L003872/1), the European Commission, H2020 Marie Skłodowska-Curie Actions
45 (HiFreq, grant no. 734317), the U.S. Department of Energy (Pacific Northwest National Lab and
46 DE-SC0019377), the National Science Foundation (grant nos. DEB-1440409, EAR-1652293,
47 EAR-1417603, and EAR-1446328), the University of Birmingham (Institute of Advanced
48 Studies), and with resources from the home institutions of the authors. Data and facilities were
49 provided by the H. J. Andrews Experimental Forest and Long Term Ecological Research
50 program, administered cooperatively by the USDA Forest Service Pacific Northwest Research
51 Station, Oregon State University, and the Willamette National Forest. In lieu of detailed author
52 contributions, we report that this study was conceptualized approximately 10 years ago and has
53 benefited tremendously from discussions with a broad group of friends and collaborators. Work
54 on this manuscript was initiated at the slow freshwater science meeting hold in Santa Maria de
55
56
57
58
59
60

1
2
3 Palautordera (Catalonia, NE Spain). The authors of this study each made specific contributions
4 to conceptualization, data collection, analysis, and/or writing and revising the manuscript. The
5 primary data analyzed are described by Ward et al. (2019) and available in Ward (2019). Results
6 of analyses completed in this study are available in Ward (2021). The authors declare no
7 conflicts of interest. Any use of trade, firm, or product names is for descriptive purposes only
8 and does not imply endorsement by the US government. Any opinions, findings, and conclusions
9 or recommendations expressed in this material are those of the authors.
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

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>

- 1
2
3 Frissell, C. A., Liss, W. J., Warren, C. E., & Hurley, M. D. (1986). A hierarchical framework for
4 stream habitat classification: Viewing streams in a watershed context. *Environmental*
5 *Management*, 10(2), 199–214.
- 6 Giraldo, L., Palacio, C., & Aguirre, N. (2014). Temporal Variation of the Extracellular
7 Enzymatic Activity (EEA): Case of Study : Aburra-Medellín River, in the Valle de Aburra
8 in Medellin, Antioquia, Colombia. *International Journal of Environmental Protection*, 4(5),
9 58–67.
- 10
11 Graham, E. B., Crump, A. R., Kennedy, D. W., Arntzen, E., Fansler, S., Purvine, S. O., et al.
12 (2018). Multi 'omics comparison reveals metabolome biochemistry, not microbiome
13 composition or gene expression, corresponds to elevated biogeochemical function in the
14 hyporheic zone. *Science of the Total Environment*, 642, 742–753.
15 <https://doi.org/10.1016/j.scitotenv.2018.05.256>
- 16 Gregory, K. J. (2006). The human role in changing river channels. *Geomorphology*, 79(3–4),
17 172–191. <https://doi.org/10.1016/j.geomorph.2006.06.018>
- 18
19 Hampton, TB., JP Zarnetske, MA Briggs, F MahmoodPoor Dehkordy, K Singha, FD Day-Lewis,
20 JW Harvey, S Roy Chowdhury and JW Lane. (2020) Experimental shifts of hydrologic
21 residence time in a sandy urban stream sediment–water interface alter nitrate removal and
22 nitrous oxide fluxes. *Biogeochemistry* 149, 195–219. [https://doi.org/10.1007/s10533-020-](https://doi.org/10.1007/s10533-020-00674-7)
23 [00674-7](https://doi.org/10.1007/s10533-020-00674-7).
- 24
25 Hampton, TB, JP Zarnetske, MA Briggs, K Singha, JW Harvey, FD Day-Lewis, F
26 MahmoodPoor Dehkordy, and JW Lane (2019) Residence time controls the fate of nitrogen
27 in flow-through lakebed sediments. *JGR-Biogeosciences*, 124, 689– 707.
28 <https://doi.org/10.1029/2018JG004741>
- 29
30 Harman, C. J., Ward, A. S., & Ball, A. (2016). How does reach-scale stream-hyporheic transport
31 vary with discharge? Insights fromrSAS analysis of sequential tracer injections in a
32 headwater mountain stream. *Water Resources Research*, 52, 7130–7150.
33 <https://doi.org/10.1002/2016WR018832>.Received
- 34
35 Harvey, J. W., & Gooseff, M. N. (2015). River corridor science: Hydrologic exchange and
36 ecological consequences from bedforms to basins. *Water Resources Research*, 51, 6893–
37 6922. <https://doi.org/10.1002/2015WR017617>
- 38
39 Heffernan, J. B., Soranno, P. A., Angilletta, M. J., Buckley, L. B., Gruner, D. S., Keitt, T. H., et
40 al. (2014). Macrosystems ecology: understanding ecological patterns and processes at
41 continental scales. *Frontiers in Ecology and the Environment*, 12(1), 5–14.
42 <https://doi.org/10.1890/130017>
- 43
44 Hill, B. H., McCormick, F. H., Harvey, B. C., Johnson, S. L., Warren, M. L., & Elonen, C. M.
45 (2010). Microbial enzyme activity, nutrient uptake and nutrient limitation in forested
46 streams. *Freshwater Biology*, 55(5), 1005–1019. [https://doi.org/10.1111/J.1365-](https://doi.org/10.1111/J.1365-2427.2009.02337.X)
47 [2427.2009.02337.X](https://doi.org/10.1111/J.1365-2427.2009.02337.X)
- 48
49 Isaak, D. J., Peterson, E. E., Ver Hoef, J. M., Wenger, S. J., Falke, J. A., Torgersen, C. E., et al.
50 (2014). Applications of spatial statistical network models to stream data. *Wiley*
51 *Interdisciplinary Reviews: Water*, 1(3), 277–294. <https://doi.org/10.1002/wat2.1023>
- 52
53 Jefferson, A., Grant, G. E., & Lewis, S. L. (2004). A River Runs Underneath It: Geological
54 Control of Spring and Channel Systems and Management Implications, Cascade Range,
55 Oregon. In *Advancing the Fundamental Sciences Proceedings of the Forest Service:*
56 *Proceedings of the Forest Service National Earth Sciences Conference* (Vol. 1, pp. 18–22).
- 57
58 Jerolmack, D. J., & Paola, C. (2010). Shredding of environmental signals by sediment transport.

- 1
2
3 *Geophysical Research Letters*, 37(19), 1–5. <https://doi.org/10.1029/2010GL044638>
- 4 Knapp, J. L. A., González-Pinzón, R., Drummond, J. D., Larsen, L. G., Cirpka, O. A., & Harvey,
5 J. W. (2017). Tracer-based characterization of hyporheic exchange and benthic biolayers in
6 streams. *Water Resources Res*, 53, 1575–1594. <https://doi.org/10.1002/2016WR019393>
- 7 Krause, S., Hannah, D. M., Fleckenstein, J. H., Heppell, C. M., Kaeser, D. H., Pickup, R., et al.
8 (2011). Inter-disciplinary perspectives on processes in the hyporheic zone. *Ecohydrology*,
9 4(4), 481–499.
- 10 Krause, S., Lewandowski, J., Grimm, N. B., Hannah, D. M., Pinay, G., McDonald, K., et al.
11 (2017). Ecohydrological interfaces as hot spots of ecosystem processes. *Water Resources*
12 *Research*, 53(8), 6359–6376. <https://doi.org/10.1002/2016WR019516>
- 13 Langbein, W. B., & Leopold, L. B. (1966). *River meanders - theory of minimum variance*.
- 14 Lee-Cullin, J. A., Zarnetske, J. P., Ruhala, S. S., & Plont, S. (2018). Toward measuring
15 biogeochemistry within the stream-groundwater interface at the network scale: An initial
16 assessment of two spatial sampling strategies. *Limnology and Oceanography: Methods*,
17 16(11), 722–733. <https://doi.org/10.1002/lom3.10277>
- 18 Leopold, L. B., Wolman, M. G., & Miller, J. P. (1964). *Fluvial Processes in Geomorphology*.
19 Dover Publications.
- 20 Leopold, L. B., & Langbein, W. B. (1962). *The Concept of Entropy in Landscape Evolution*.
- 21 Li, L., Sullivan, P. L., Benettin, P., Cirpka, O. A., Bishop, K., Brantley, S. L., et al. (2021).
22 Toward catchment hydro-biogeochemical theories. *Wiley Interdisciplinary Reviews: Water*,
23 8(1), e1495. <https://doi.org/10.1002/wat2.1495>
- 24 Liébault, F., & Piégay, H. (2002). Causes of 20th century channel narrowing in mountain and
25 piedmont rivers of southeastern France. *Earth Surface Processes and Landforms*, 27(4),
26 425–444. <https://doi.org/10.1002/esp.328>
- 27 Lovett, G. M., Burns, D. A., Driscoll, C. T., Jenkins, J. C., Mitchell, M. J., Rustad, L., et al.
28 (2007). Who needs environmental monitoring? *Frontiers in Ecology and the Environment*,
29 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)
- 30 Martin, P. Y., & Turner, B. A. (1986). Grounded Theory and Organizational Research. *The*
31 *Journal of Applied Behavioral Science*, 22(2), 141–157.
32 <https://doi.org/10.1177/002188638602200207>
- 33 McCluney, K. E., Poff, N. L., Palmer, M. A., Thorp, J. H., Poole, G. C., Williams, B. S., et al.
34 (2014). Riverine macrosystems ecology: sensitivity, resistance, and resilience of whole river
35 basins with human alterations. *Frontiers in Ecology and the Environment*, 12(1), 48–58.
36 <https://doi.org/10.1890/120367>
- 37 McGuire, K. J., Torgersen, C. E., Likens, G. E., Buso, D. C., Lowe, W. H., & Bailey, S. W.
38 (2014). Network analysis reveals multiscale controls on streamwater chemistry.
39 *Proceedings of the National Academy of Sciences of the United States of America*, 111(19),
40 7030–7035. <https://doi.org/10.1073/pnas.1404820111>
- 41 Minshall, G. W., Petersen, R. C., Cummins, K. W., Bott, T. L., Sedell, J. R., Cushing, C. E., &
42 Vannote, R. L. (1983). Interbiome Comparison of Stream Ecosystem Dynamics. *Ecological*
43 *Monographs*, 53(1), 1–25. <https://doi.org/10.2307/1942585>
- 44 Nikora, V. I., & Hicks, D. M. (1997). Scaling Relationships for Sand Wave Development in
45 Unidirectional Flow. *Journal of Hydraulic Engineering*, 123(12), 1152–1156.
46 [https://doi.org/10.1061/\(asce\)0733-9429\(1997\)123:12\(1152\)](https://doi.org/10.1061/(asce)0733-9429(1997)123:12(1152))
- 47 O'Malley, M.A. (2008). “Everything is everywhere: but the environment selects”: ubiquitous
48 distribution and ecological determinism in microbial biogeography. *Studies in History and*
49
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3 *Philosophy of Biological and Biomedical Sciences*, 39(3), 314–325.
4 <https://doi.org/10.1016/J.SHPSC.2008.06.005>
5
6 Payn, R. A., Gooseff, M. N., McGlynn, B. L., Bencala, K. E., & Wondzell, S. M. (2009).
7 Channel water balance and exchange with subsurface flow along a mountain headwater
8 stream in Montana, United States. *Water Resources Research*, 45.
9
10 Pinay, G., Peiffer, S., De Dreuzy, J. R., Krause, S., Hannah, D. M., Fleckenstein, J. H., et al.
11 (2015). Upscaling Nitrogen Removal Capacity from Local Hotspots to Low Stream Orders'
12 Drainage Basins. *Ecosystems*, 18(6), 1101–1120. [https://doi.org/10.1007/s10021-015-9878-](https://doi.org/10.1007/s10021-015-9878-5)
13 5
14 Pringle, C. M., Naiman, R. J., Bretschko, G., Karr, J. R., Oswood, M. W., Webster, J. R., et al.
15 (1988). Patch Dynamics in Lotic Systems: The Stream as a Mosaic. *Journal of the North*
16 *American Benthological Society*, 7(4), 503–524. <https://doi.org/10.2307/1467303>
17
18 Rana, S. M. M., Scott, D. T., & Hester, E. T. (2017). Effects of in-stream structures and channel
19 flow rate variation on transient storage. *Journal of Hydrology*, 548, 157–169.
20 <https://doi.org/10.1016/j.jhydrol.2017.02.049>
21
22 Robertson, A. D., Paustian, K., Ogle, S., Wallenstein, M. D., Lugato, E., & Francesca Cotrufo,
23 M. (2019). Unifying soil organic matter formation and persistence frameworks: The MEMS
24 model. *Biogeosciences*, 16(6), 1225–1248. <https://doi.org/10.5194/BG-16-1225-2019>
25
26 Rodríguez-Iturbe, I., & Rinaldo, A. (1997). *Fractal River Basins: Chance and Self-Organization*.
27 Cambridge, UK: Cambridge University Press.
28
29 Santschi, P. H., Presley, B. J., Wade, T. L., Garcia-Romero, B., & Baskaran, M. (2001).
30 Historical contamination of PAHs, PCBs, DDTs, and heavy metals in Mississippi River
31 Delta, Galveston Bay and Tampa Bay sediment cores. *Marine Environmental Research*,
32 52(1), 51–79. [https://doi.org/10.1016/S0141-1136\(00\)00260-9](https://doi.org/10.1016/S0141-1136(00)00260-9)
33
34 Sinsabaugh, R. L., Findlay, S., Franchini, P., & Fischer, D. (1997). Enzymatic analysis of
35 riverine bacterioplankton production. *Limnology and Oceanography*, 42(1), 29–38.
36 <https://doi.org/10.4319/LO.1997.42.1.0029>
37
38 Sinsabaugh, R. L., Findlay, S., Franchini, P., & Fischer, D. (1997). Enzymatic analysis of
39 riverine bacterioplankton production. *Limnology and Oceanography*, 42(1), 29–38.
40 <https://doi.org/10.4319/LO.1997.42.1.0029>
41
42 Sinsabaugh, R. L., & Shah, J. J. F. (2012). Ecoenzymatic Stoichiometry and Ecological Theory.
43 <http://Dx.Doi.Org/10.1146/Annurev-Ecolsys-071112-124414>, 43, 313–343.
44 <https://doi.org/10.1146/ANNUREV-ECOLSYS-071112-124414>
45
46 Smidt, S. J., Cullin, J. A., Ward, A. S., Robinson, J., Zimmer, M. A., Lutz, L. K., & Endreny, T.
47 A. (2015). A Comparison of Hyporheic Transport at a Cross-Vane Structure and Natural
48 Riffle. *Ground Water*, 53(6), 859–871. <https://doi.org/10.1111/gwat.12288>
49
50 Sollins, P., Cromack, K., Corison, F. M. M., Waring, R. H., & Harr, R. D. (1981). Changes in
51 Nitrogen Cycling at an Old-Growth Douglas-fir Site After Disturbance. *Journal of*
52 *Environmental Quality*, 10(1), 37–42.
53 <https://doi.org/10.2134/JEQ1981.00472425001000010007X>
54
55 Stegen, J. C., & Goldman, A. E. (2018). WHONDRS: a Community Resource for Studying
56 Dynamic River Corridors. *MSystems*, 3(5), 151–169.
57 <https://doi.org/10.1128/msystems.00151-18>
58
59 Stegen, J. C., Bottos, E. M., & Jansson, J. K. (2018). A unified conceptual framework for
60 prediction and control of microbiomes. *Current Opinion in Microbiology*, 44(July), 20–27.
<https://doi.org/10.1016/j.mib.2018.06.002>

- 1
2
3 Stegen, J. C., Johnson, T., Fredrickson, J. K., Wilkins, M. J., Konopka, A. E., Nelson, W. C., et
4 al. (2018). Influences of organic carbon speciation on hyporheic corridor biogeochemistry
5 and microbial ecology. *Nature Communications*, 9(1), 1–11.
6 <https://doi.org/10.1038/s41467-018-03572-7>
- 8 Strauss, A., & Corbin, J. (1994). Grounded theory methodology: An overview. In N. Denzin &
9 Y. Lincoln (Eds.), *Handbook of qualitative research* (pp. 273–285). Sage Publications, Inc.
- 11 Swanson, F. J., & James, M. E. (1975). *Geology and geomorphology of the H.J. Andrews*
12 *Experimental Forest, western Cascades, Oregon*. Portland, OR.
- 13 Swanson, F. J., & Jones, J. A. (2002). Geomorphology and hydrology of the H.J. Andrews
14 Experimental Forest, Blue River, Oregon. In *Field guide to geologic processes in Cascadia*.
- 15 Tank, J. L., Rosi-Marshall, E. J., Baker, M. A., & Hall, R. O. (2008). Are rivers just big streams?
16 A pulse method to quantify nitrogen demand in a large river. *Ecology*, 89(10), 2935–2945.
- 17 Tiegs, S. D., Costello, D. M., Isken, M. W., Woodward, G., McIntyre, P. B., Gessner, M. O., et
18 al. (2019). Global patterns and drivers of ecosystem functioning in rivers and riparian zones.
19 *Science Advances*, 5(1), eaav0486. <https://doi.org/10.1126/SCIADV.AAV0486>
- 20 Triska, F. J., Sedell, J. R., Cromack, K., Gregory, S. V., & McCorison, F. M. (1984). Nitrogen
21 Budget for a Small Coniferous Forest Stream. *Ecological Monographs*, 54(1), 119–140.
22 <https://doi.org/10.2307/1942458>
- 23 Turnbull, L., Hütt, M. T., Ioannides, A. A., Kininmonth, S., Poepl, R., Tockner, K., et al. (2018,
24 December 1). Connectivity and complex systems: learning from a multi-disciplinary
25 perspective. *Applied Network Science*. Springer. <https://doi.org/10.1007/s41109-018-0067-2>
- 26 Valett, H. M., Morrice, J. A., Dahm, C. N., & Campana, M. E. (1996). Parent lithology, surface-
27 groundwater exchange, and nitrate retention in headwater streams. *Limnology and*
28 *Oceanography*, 333–345.
- 29 Vannote, R. L., Minshall, G. W., Cummins, K. W., Sedell, J. R., & Cushing, C. E. (1980). The
30 River Continuum Concept. *Canadian Journal of Fisheries and Aquatic Sciences*, 37, 130–
31 137.
- 32 Ver Hoef, J. M., Peterson, E., & Theobald, D. (2006). Spatial statistical models that use flow and
33 stream distance. *Environmental and Ecological Statistics*, 13(4), 449–464.
34 <https://doi.org/10.1007/s10651-006-0022-8>
- 35 Walling, D. E., & Fang, D. (2003). Recent trends in the suspended sediment loads of the world's
36 rivers. *Global and Planetary Change*, 39(1–2), 111–126. [https://doi.org/10.1016/S0921-](https://doi.org/10.1016/S0921-8181(03)00020-1)
37 [8181\(03\)00020-1](https://doi.org/10.1016/S0921-8181(03)00020-1)
- 38 Wallis, I., Prommer, H., Berg, M., Siade, A. J., Sun, J., & Kipfer, R. (2020). The river–
39 groundwater interface as a hotspot for arsenic release. *Nature Geoscience*, 13(4), 288–295.
40 <https://doi.org/10.1038/s41561-020-0557-6>
- 41 Ward, A. S. (2015). The evolution and state of interdisciplinary hyporheic research. *Wiley*
42 *Interdisciplinary Reviews: Water*, 3(1), 83–103. <https://doi.org/10.1002/wat2.1120>
- 43 Ward, A. S. (2019). ESSD, 2019 - Data Collection. <https://doi.org/10.5194/essd-11-1-2019>
- 44 Ward, A. S. (2021). Supporting data for Ward et al., (In Review) Advancing river corridor
45 science beyond disciplinary boundaries with an inductive approach to hypothesis
46 generation, HydroShare, Accessed 6-May-2021.
47 <http://www.hydroshare.org/resource/de6d92d314354ea6819157818669fc59>
- 48 Ward, A. S., & Packman, A. I. (2019). Advancing our predictive understanding of river corridor
49 exchange. *Wiley Interdisciplinary Reviews: Water*, 6(1), e1327.
50 <https://doi.org/10.1002/wat2.1327>
- 51
52
53
54
55
56
57
58
59
60

- 1
2
3 Ward, A. S., Zarnetske, J. P., Baranov, V., Blaen, P. J., Brekenfeld, N., Chu, R., et al. (2019).
4 Co-located contemporaneous mapping of morphological, hydrological, chemical, and
5 biological conditions in a 5th-order mountain stream network, Oregon, USA. *Earth System*
6 *Science Data*, 11(4). <https://doi.org/10.5194/essd-11-1567-2019>
7
8 Ward, A. S., Wondzell, S. M., Schmadel, N. M., Herzog, S., Zarnetske, J. P., Baranov, V., et al.
9 (2019). Spatial and temporal variation in river corridor exchange across a 5th order
10 mountain stream network. *Hydrology and Earth System Sciences Discussions*, (April), 1–
11 39. <https://doi.org/10.5194/hess-2019-108>
12
13 Waring, B. G., Sulman, B. N., Reed, S., Smith, A. P., Averill, C., Creamer, C. A., et al. (2020).
14 From pools to flow: The PROMISE framework for new insights on soil carbon cycling in a
15 changing world. *Global Change Biology*, 26(12), 6631–6643.
16 <https://doi.org/10.1111/GCB.15365>
17
18 Williams, C. J., Scott, A. B., Wilson, H. F., & Xenopoulos, M. A. (2011). Effects of land use on
19 water column bacterial activity and enzyme stoichiometry in stream ecosystems. *Aquatic*
20 *Sciences 2011* 74:3, 74(3), 483–494. <https://doi.org/10.1007/S00027-011-0242-3>
21
22 Williams, C. J., Yamashita, Y., Wilson, H. F., Jaffé, R., & Xenopoulos, M. A. (2010).
23 Unraveling the role of land use and microbial activity in shaping dissolved organic matter
24 characteristics in stream ecosystems. *Limnology and Oceanography*, 55(3), 1159–1171.
25 <https://doi.org/10.4319/LO.2010.55.3.1159>
26
27 Wisnoski, N. I., & Lennon, J. T. (2021). Microbial community assembly in a multi-layer
28 dendritic metacommunity. *Oecologia*, 195(1), 13–24. [https://doi.org/10.1007/s00442-020-](https://doi.org/10.1007/s00442-020-04767-w)
29 [04767-w](https://doi.org/10.1007/s00442-020-04767-w)
30
31 Wisnoski, N. I., Muscarella, M. E., Larsen, M. L., Peralta, A. L., & Lennon, J. T. (2020).
32 Metabolic insight into bacterial community assembly across ecosystem boundaries.
33 *Ecology*, 101(4), e02968. <https://doi.org/10.1002/ECY.2968>
34
35 Wit, R. De, & Bouvier, T. (2006). ‘Everything is everywhere, but, the environment selects’; what
36 did Baas Becking and Beijerinck really say? *Environmental Microbiology*, 8(4), 755–758.
37 <https://doi.org/10.1111/J.1462-2920.2006.01017.X>
38
39 Wohl, E. (2005). Compromised Rivers: Understanding Historical Human Impacts on Rivers in
40 the Context of Restoration. *Ecology and Society*, 10(2), 2.
41
42 Wondzell, S. M., & Gooseff, M. N. (2014). Geomorphic Controls on Hyporheic Exchange
43 Across Scales: Watersheds to Particles. In J. Schroder & E. Wohl (Eds.), *Treatise on*
44 *Geomorphology* (Vol. 9, pp. 203–218). San Diego, CA: Academic Press.
45
46 Wood, P. J., Boulton, A. J., Little, S., & Stubbington, R. (2010). Is the hyporheic zone a
47 refugium for aquatic macroinvertebrates during severe low flow conditions? *Fundamental*
48 *and Applied Limnology / Archiv Für Hydrobiologie*, 176(4), 377–390.
49 <https://doi.org/10.1127/1863-9135/2010/0176-0377>
50
51 Wörman, A., Packman, A. I., Marklund, L., Harvey, J. W., & Stone, S. H. (2007). Fractal
52 topography and subsurface water flows from fluvial bedforms to the continental shield.
53 *Geophysical Research Letters*, 34(7), 1–5. <https://doi.org/10.1029/2007GL029426>
54
55 Wu, L., Singh, T., Gomez-Velez, J., Nützmann, G., Wörman, A., Krause, S., & Lewandowski, J.
56 (2018). Impact of Dynamically Changing Discharge on Hyporheic Exchange Processes
57 Under Gaining and Losing Groundwater Conditions. *Water Resources Research*, 54(12),
58 10,076–10,093. <https://doi.org/10.1029/2018WR023185>
59
60 Yoder, L., Ward, A. S., Spak, S., & Dalrymple, K. (2020). Local Government Perspectives on
Collaborative Governance: A Comparative Analysis of Iowa’s Watershed Management

1
2
3 Authorities. *Policy Studies Journal*. <https://doi.org/10.1111/psj.12389>
4 Zhou, C., Liu, Y., Liu, C., Liu, Y., & Tfaily, M. M. (2019). Compositional changes of dissolved
5 organic carbon during its dynamic desorption from hyporheic zone sediments. *Science of*
6 *the Total Environment*, 658, 16–23. <https://doi.org/10.1016/j.scitotenv.2018.12.189>
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

Supplemental Material for

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

Adam Ward¹, Aaron Packman², Susana Bernal³, Nicolai Brekenfeld⁴, Jen Drummond⁴, Emily Graham⁵, David M. Hannah⁴, Megan Klaar⁶, Stefan Krause⁴, Marie Kurz⁷, Angang Li², Anna Lupon³, Feng Mao⁸, M. Eugènia Martí Roca³, Valerie Ouellet⁴, Todd Royer¹, James Stegen⁵, Jay Zarnetske⁹

1 O'Neill School of Public and Environmental Affairs, Indiana University, Bloomington, Indiana, USA

2 Department of Civil and Environmental Engineering, Northwestern University, Evanston, Illinois, USA

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

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

5 Earth and Biological Sciences Division, Pacific Northwest National Laboratory, Richland, Washington, USA

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

7 The Academy of Natural Sciences of Drexel University, Philadelphia, Pennsylvania, USA

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

9 Department of Earth and Environmental Sciences, Michigan State University, East Lansing, Michigan, USA

Contents of this file

Figures S1 to S3

Tables S1 to S3

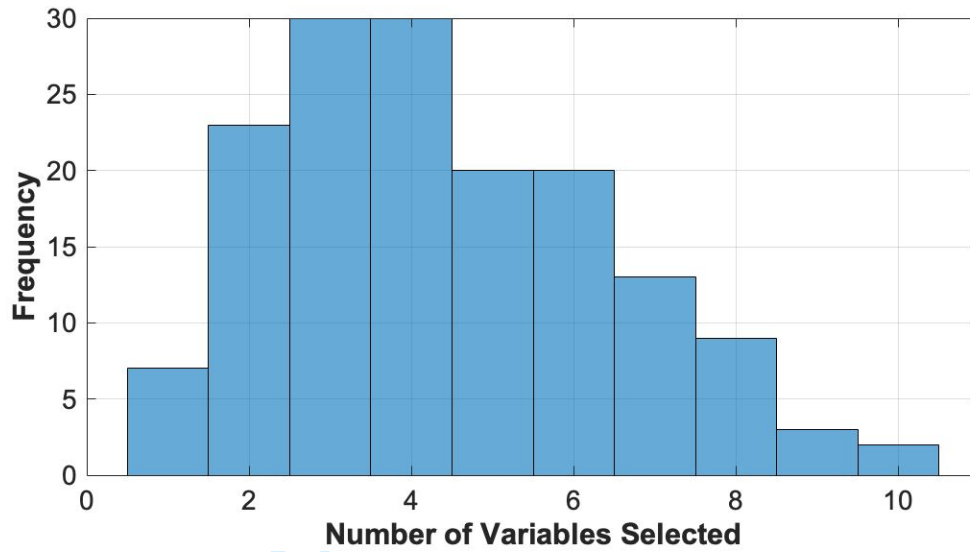


Figure S1. Histogram of the number of variables selected for each SVMR constructed on all variables, demonstrating models were parsimonious.

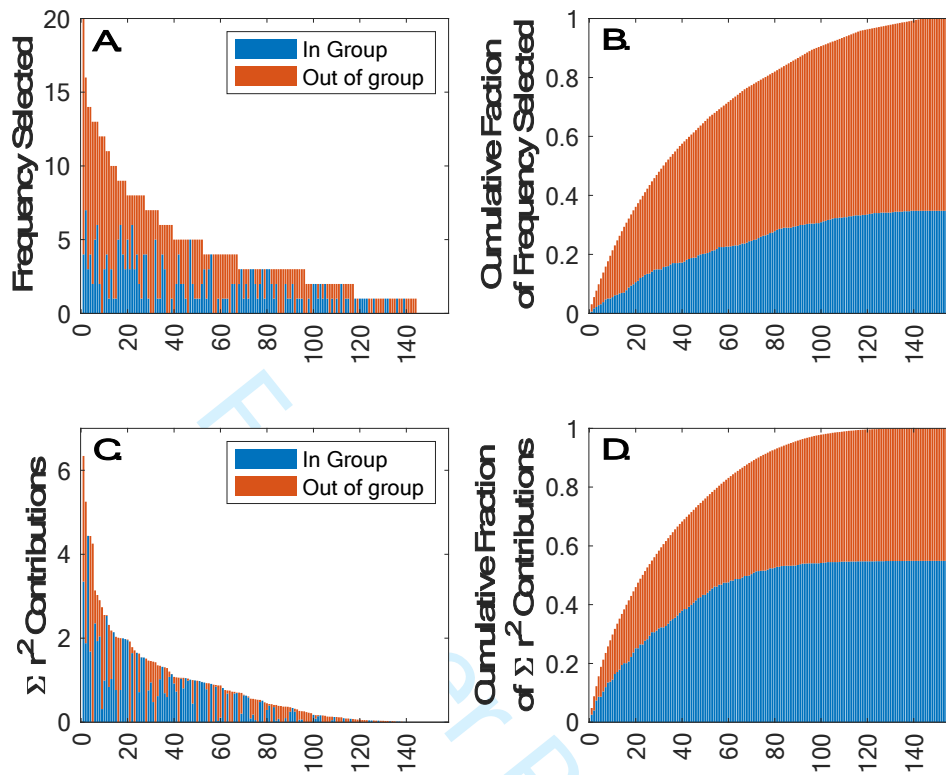


Figure S2. In- and out-of-group selection frequency (top row) and contributions to model r^2 (bottom row) for SVMRs constructed on all variables. Variables are ordered from most- to least-frequently selected on the x-axis. Data are displayed for individual variables (left column) and cumulatively (right column).

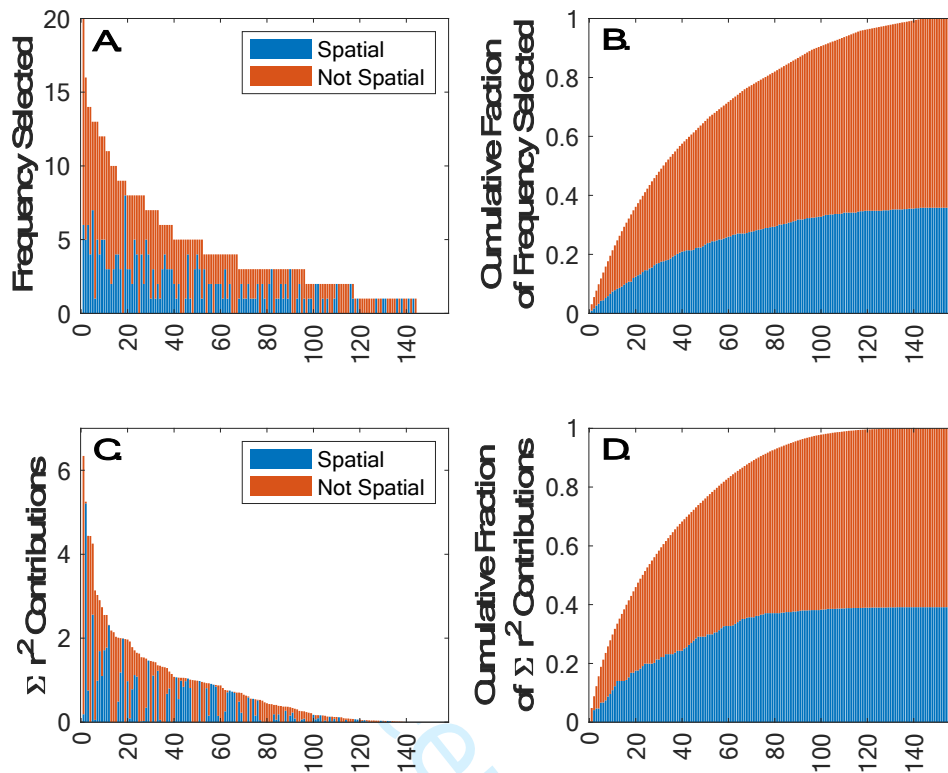


Figure S3. Spatially structured and not structured selection frequency (top row) and contributions to model r^2 (bottom row) for SVMs constructed on all variables. Variables are ordered from most- to least-frequently selected on the x-axis. Data are displayed for individual variables (left column) and cumulatively (right column).

	DOM-S	GEO-S	ICR-S	MACRO-S	NUTS-S	PCHEM-S	TRACER-S	DOM-NOT	GEO-NOT	ICR-NOT	MACRO-NOT	NUTS-NOT	PCHEM-NOT	TRACER-NOT
DOM-S	0	0	0	0	1	0	2	1	4	0	2	1	0	0
GEO-S	0	24	0	0	1	2	10	1	31	2	15	0	0	2
ICR-S	0	0	5	0	0	0	0	0	3	3	0	0	0	2
MACRO-S	0	0	0	0	0	0	1	0	0	0	5	0	0	1
NUTS-S	0	5	1	1	3	0	9	2	18	0	6	2	0	3
PCHEM-S	0	1	0	0	0	1	6	0	14	1	3	0	0	1
TRACER-S	0	5	0	1	2	1	6	1	23	1	6	0	0	3
DOM-NOT	1	1	0	1	1	0	1	4	9	3	7	1	1	2
GEO-NOT	0	23	0	5	2	0	12	5	68	2	28	1	2	14
ICR-NOT	0	1	1	3	0	0	2	0	23	13	7	0	0	7
MACRO-NOT	0	6	1	2	2	1	5	2	27	3	35	0	0	7
NUTS-NOT	0	1	0	0	1	0	2	2	8	2	4	1	0	2
PCHEM-NOT	0	1	0	1	0	0	1	0	11	1	2	0	3	2
TRACER-NOT	0	3	0	1	0	1	7	1	28	4	6	0	0	10

Table S1. Summary of the frequency with which variables from one disciplinary subgroup (column) were used as a predictor for variables in each disciplinary subgroups (rows). Yellow diagonals represent within-subgroup information flows. This table is the tabular representation of relationships shown in Fig. 2 of the manuscript.

		Variable being predicted						
		DOM	GEO	ICR	MACRO	Nutrients	PChem	tracer
Source of information	DOM-PC1	10	45	0	6	1	3	6
	DOM-PC2	2	45	4	8	0	0	9
	GEO-PC1	0	59	0	0	0	7	8
	GEO-PC2	6	57	4	13	5	0	3
	ICR-PC1	0	28	14	9	0	3	8
	ICR-PC2	3	8	15	6	6	0	3
	MACRO-PC1	0	20	5	19	10	0	0
	MACRO-PC2	0	26	1	23	6	0	1
	Nutrients-PC1	0	31	3	2	6	0	9
	Nutrients-PC2	2	24	6	5	5	0	7
	PChem-PC1	2	49	0	0	5	12	15
	PChem-PC2	0	33	5	0	5	8	11
	tracer-PC1	0	47	4	5	0	4	18
	tracer-PC2	5	61	1	4	1	7	25
No. of PCs Used?		7	14	11	11	10	7	13
No. of Groups?		6	7	7	6	7	5	7

Table S2. Summary of the frequency with which PC1 and PC2 from each disciplinary subgroup (rows) were used as a predictor for variables in each disciplinary subgroups (columns). Yellow diagonals represent within-subgroup information flows. This table is the tabular representation of relationships shown in Fig. 4 of the manuscript.

	DOM	GEO	ICR	MACRO	NUTRIENT	PCHEM	TRACER
DOM	190 (<0.1%) 3 (0.48%)	1594 (<0.1%) 20 (3.2%)	241 (<0.1%) 3 (0.47%)	0 (0%) 12 (1.9%)	46 (<0.1%) 8 (1.3%)	29 (<0.1%) 1 (0.15%)	199 (<0.1%) 7 (1.11%)
GEO		715257 (26.2%) 125 (19.9%)	15031 (0.55%) 30 (4.8%)	1325 (<0.1%) 80 (12.7%)	54621 (2.0%) 36 (5.7%)	27451 (1.0%) 30 (4.8%)	260085 (9.5%) 95 (15.1%)
ICR			7385 (0.27%) 15 (2.4%)	176 (<0.1%) 13 (2.1%)	17441 (0.64%) 3 (0.48%)	10995 (0.40%) 2 (0.32%)	22239 (0.81%) 16 (2.5%)
MACRO				650 (<0.1%) 27 (4.3%)	639 (<0.1%) 13 (2.1%)	203 (<0.1%) 7 (1.1%)	1154 (<0.1%) 27 (4.3%)
NUTRIENT					57790 (2.1%) 6 (0.95%)	86107 (3.2%) 0 (0%)	552278 (20.2%) 18 (2.9%)
PCHEM						129536 (4.7%) 2 (0.32%)	117384 (4.3%) 12 (1.9%)
TRACER							651648 (23.9%) 18 (2.9%)

Table S3. Frequency of pairwise study in literature assessment and SVMRs, organized by group. Values are listed as the raw number and percentage of each search each cell.