

# Physical habitat is more than a sediment issue: A multi-dimensional habitat assessment indicates new approaches for river management

#### STAC Workshop: Leveraging AI and ML 24 February 2025

#### Matthew Cashman

US Geological Survey Water Mission Area Earth Surface Processes Division mcashman@usgs.gov



## Why ML? Predicting habitat quality in unmonitored areas





#### Research article

Physical habitat is more than a sediment issue: A multi-dimensional habitat assessment indicates new approaches for river management

Matthew J. Cashman<sup>a,\*</sup>, Gina Lee<sup>b</sup>, Leah E. Staub<sup>b</sup>, Michelle P. Katoski<sup>c</sup>, Kelly O. Maloney<sup>d</sup>

Study goals:

- 1. Where is physical habitat good/bad?
  - Supervised continuous predictions
- 2. Are there distinct dimensions of physical habitat?
  - Unsupervised dimensionality reduction and clustering
- 3. How is habitat affected by available management intervention pathways?





## Why ML? Predictions allow model overlays





Is degraded habitat caused by a sediment supply problem?

- Will restricting sediment supply improve habitat?
- What about confounding problems with flow alteration?

Restorations that focus on restricting sediment, without addressing flows or in-channel hydromorphic diversity, are unlikely to improve the habitat metrics that justified the TMDL.



#### Why ML? They are only the means, not the end



Focus of paper is management implications, not ML methods

> If ML methods distracted from that focus, they were omitted

Tree-based: **random forest**, XGBoost, lightGBM, H2O AutoML, and ensemble stacking ~*comparable performances* 

- >20,000 observations, synoptic design
- Adjusted for mean-centered bias (Belitz and Stackelberg, 2021)

Tested many explainable AI techniques, most not in paper

- · Expected relationships, nothing 'novel'
- Used for internal model consistency/validation





Preliminary Information-Subject to Revision. Not for Citation or Distribution

# Aligning with CBP goals? Restoration targeting and design

Modeled metrics are routinely used in field monitoring and well-known by stakeholders.

Models can help prioritize and identify areas for restoration or conservation (by itself or with co-occurring stressors).





## Challenges: Input data quality

Limited by quality of input data

- Metrics are visually scored, semi-quantitative, ٠ subjective
- Field-measurement uncertainty accounted for • ~80% of RMSE in our ML models



http://www.epa.gov/OWOW/monitoring/techmon.html

By:	Project Officer:					
Michael T. Barbour	Chris Faulkner					
Jeroen Gerritsen	Office of Water					
Blaine D. Snyder	USEPA					
James B. Stribling	401 M Street, NW					

2a. Embeddedness-High Gradient





**Optimal** Range

(William Taft, MI DNR) Poor Range

(William Taft, MI DNR)

	Habitat	Condition Category																			
	Parameter	Optimal					Su	bopti		Marginal					Poor						
	2.a Embeddedness (high gradient)	Gravel, cobble, and boulder particles are 0- 25% surrounded by fine sediment. Layering of cobble provides diversity of				Gravel, cobble, and boulder particles are 25- 50% surrounded by fine sediment.					Gravel, cobble, and boulder particles are 50- 75% surrounded by fine sediment.				Gravel, cobble, and boulder particles are more than 75% surrounded by fine sediment.						
l	SCORE	20 1	9 18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0



6

#### Challenges: Causal questions need causal methods

#### Judea Pearl's Causal Hierarchy

	Layer (Symbolic)	Typical Activity	Typical Question	Example Research Question	Statistical Methods	Modified from Bareinbom et al 2022 <sup>1</sup>
L <sub>1</sub>	Associational $P(y x)$	Seeing	What is?	Where are areas quality physical habitat across the Chesapeake Bay watershed?	Supervised / Unsupervised ML	Traditional non-causal ML
L <sub>2</sub>	Interventional $P(y do(x),c)$	Doing	What if I do?	How does dam removal affect physical habitat?	Reinforcement Learning Randomized Controlled Trials A/B Testing "Observational Experiments"	Caucal
L <sub>3</sub>	Counterfactual $P(y_x x',y')$	Imagining	What if instead? Why?	What would physical habitat be if there was no anthropogenic effects at all?	Causal Mediation/Path Synthetic Control <b>Causal ML</b>	methods

- Questions at higher layers cannot be accurately answered with information and methods from lower levels
- Traditional ML can be accurate for "What is?" Qs  $(L_1)$  and inaccurate for "What if?" or counterfactual scenario Qs  $(L_2, L_3)$ 
  - ML methods can infer info for a variable without it being in the dataset (Kratzert et al., 2019<sup>2</sup>)
  - Correlation, confounding, and hierarchical dependency causes problems for estimating cause-effect L<sub>2</sub>, L<sub>3</sub> scenarios, less so for L<sub>1</sub> predictions
- But what about progress-guided DL models pre-trained on cause-effect process models?



## Challenges: Causal machine learning is rapidly developing

- Causal Inference, Causal Machine Learning, and Causal Discovery
  - Used largely in public health, econometrics, genomics, nascent in ecology (some private co., Google, Uber, Microsoft)
- Causal inference/ML techniques are specifically designed to accurately estimate cause and effect
  - Propensity score matching/weighting
  - Causal forests
  - Double-machine learning
  - Targeted Maximum Likelihood Estimation (TMLE)
  - Highly Adaptive LASSO (HAL)
  - Causal Impact (using Bayesian structural time-series models)
  - Deep End-to-end Causal Inference (DECI)
  - Among many others...
- This is a rapidly developing field, and not all methods are suitable (yet)
  - Depends on your dataset, specific questions, and assumptions
  - · Lots of nuance, little "off-the-shelf" accessibility



Useful intro texts to causal inference



Scott Cunningham