

# *Machine learning and differentiable modeling for Geosciences*

**Chaopeng Shen**

<sup>1</sup>Civil and Environmental Engineering  
**Penn State University**  
**cshen@engr.psu.edu**



<https://github.com/mhpi>  
Hydroml.org

Hydroml.org  
HydroML Symposium, May 22-26, 2022, Penn State  
HydroML 2, May 2023, Berkeley, CA

# About me

- Ph.D. Michigan State in Env. Engr.
- Postdoc Lawrence Berkeley National Lab
- Associate Editor, Water Resources Research  
Specialty Chief Editor, Frontiers in Water:  
Water and AI.
- “Grew up” as a process-based modeler,  
solving PDEs. See both sides of the story.
- Got into ML since 2016.





# Overview

- **What** is the fundamental strengths of ML models compared to process-based models?
- **What** is *differentiable modeling (DM) in geosciences*?
- **What** can DM bring into global hydrology?


*Shen et al., 2023 Nature Reviews Earth & Environment* <https://t.co/qyuAzYPA6Y>

---

nature reviews earth & environment

<https://doi.org/10.1038/s43017-023-00450-9>

Perspective

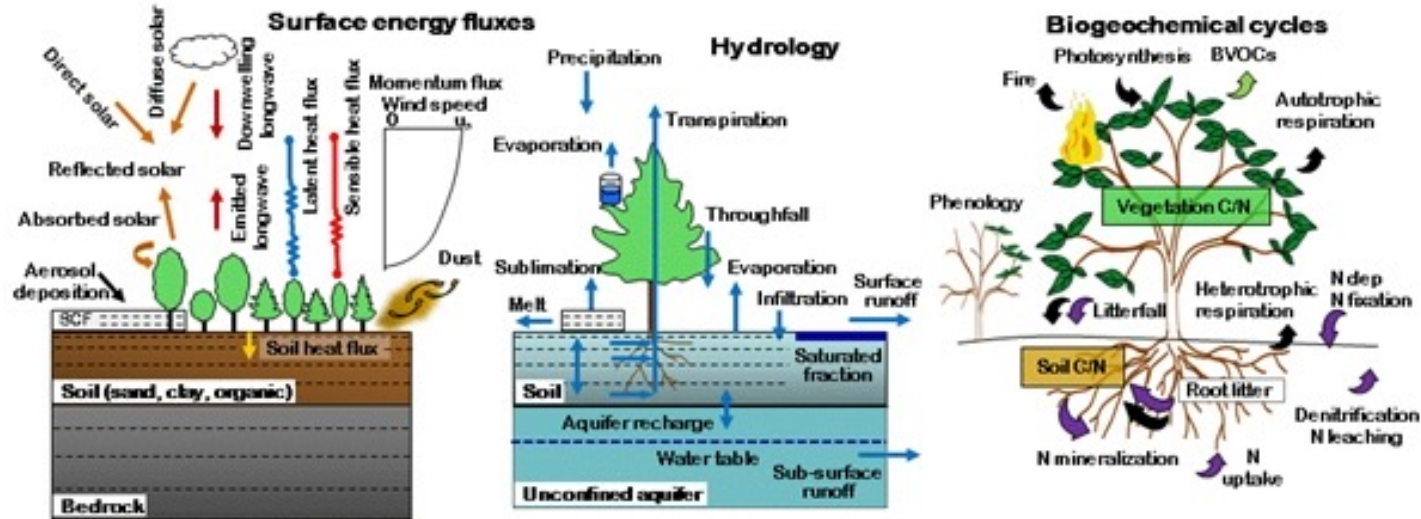
 Check for updates

---

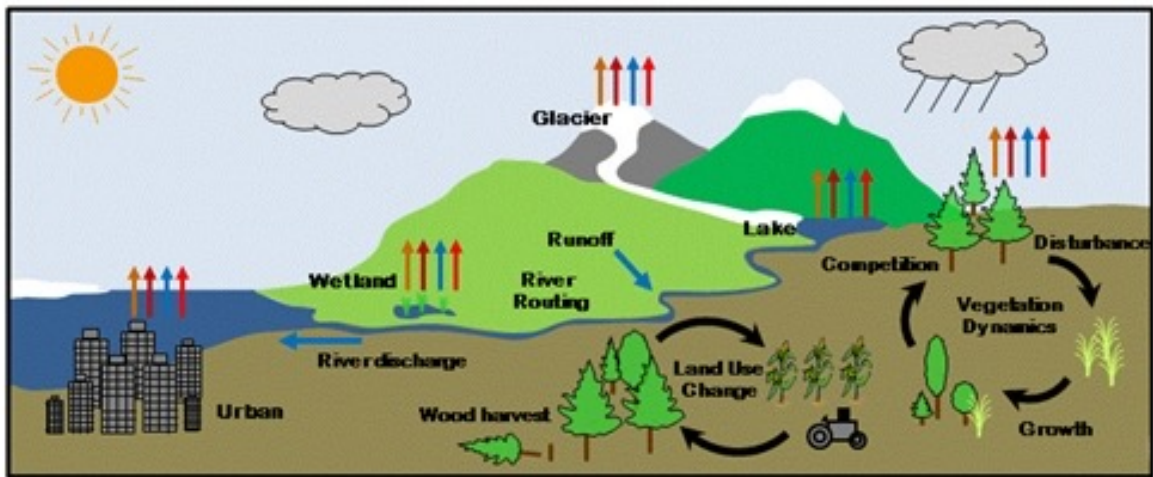
## Differentiable modelling to unify machine learning and physical models for geosciences

A list of authors and their affiliations appears at the end of the paper

# Process-based Earth-system models were highly valuable but some challenges emerged...



- Increasing complexity
- Difficult to evolve quickly with more big data.
- May contain problematic assumptions.
- Influenced by human intuition & biases





# What is DL and why DL?

a rebranding of neural networks featuring

- (i) Large capacity
- (ii) Hidden layers that automatically extract features
- (iii) Improved architecture/regularization
- (iv) Working directly with data

a primary value proposition is the avoidance of expertise!

## Three phases

1. Use ML to learn where the limit is.
2. Understand the gaps in our knowledge.
3. Using ML to unify across domains.



## Water Resources Research

AN AGU JOURNAL

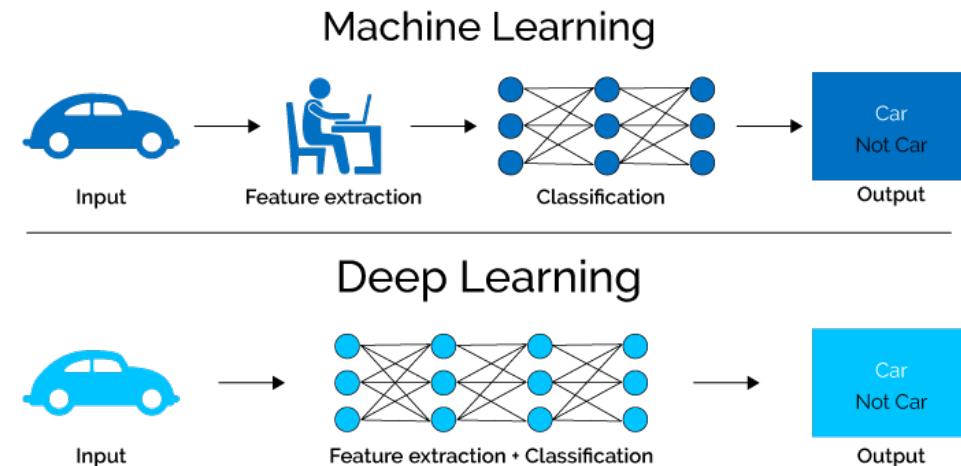
Review Article | [Open Access](#)

A trans-disciplinary review of deep learning research and its relevance for water resources scientists

Chaopeng Shen [✉](#)

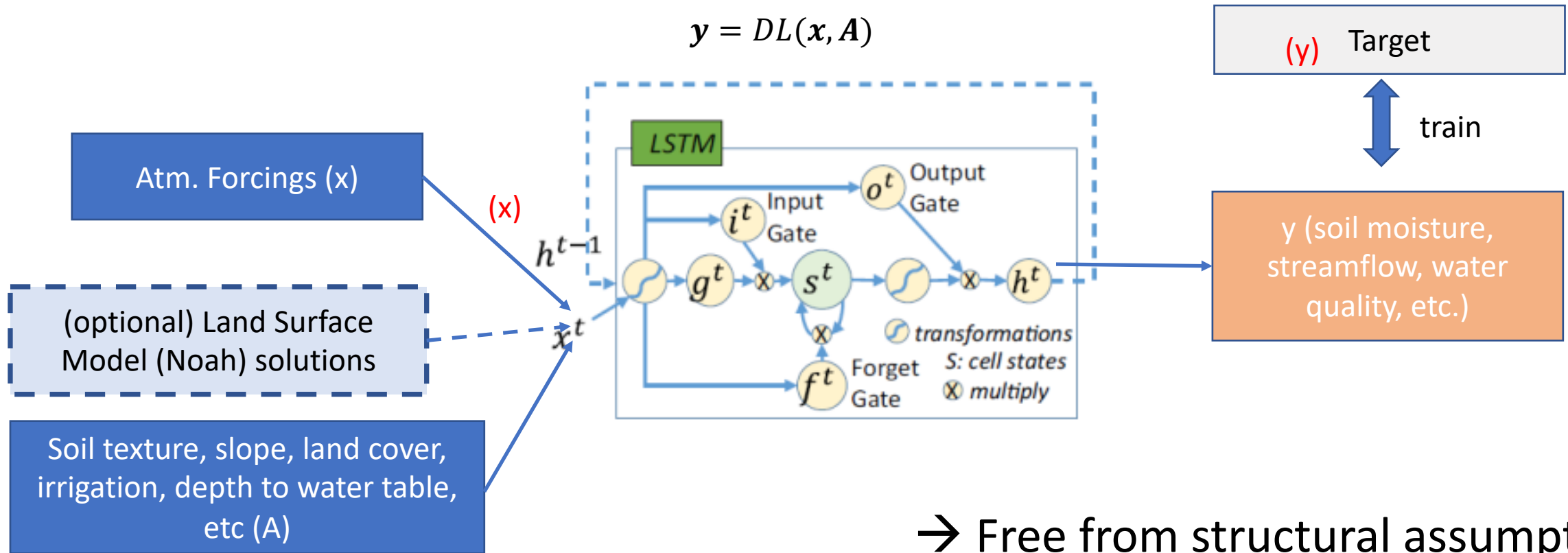
First published: 30 August 2018 | <https://doi.org/10.1029/2018WR022643>

$X \rightarrow Y$



# Hydrologic DL phase 1.

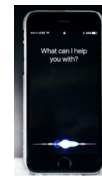
A hydrologic model w/o structural assumptions...



→ Free from structural assumptions

→ A chance to start anew!

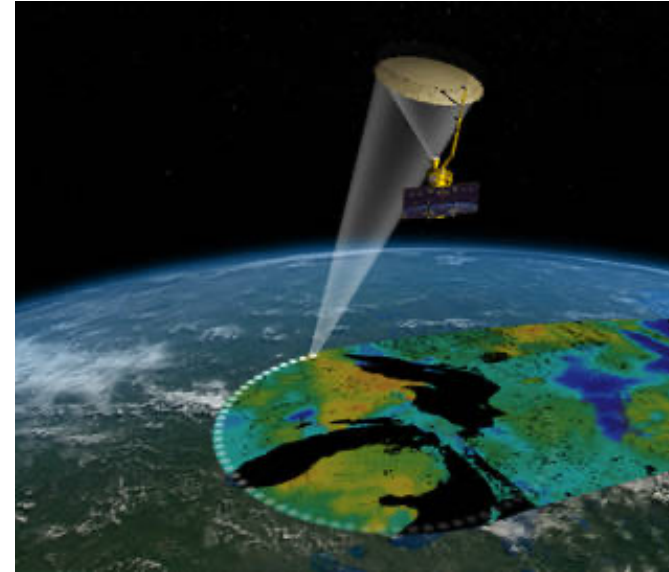
→ A chance to see where the limit is!



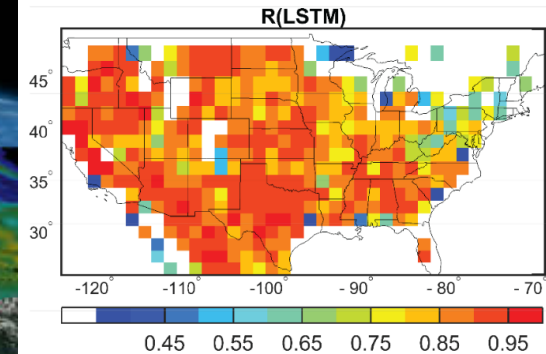


# Case studies– first phase of DL in water

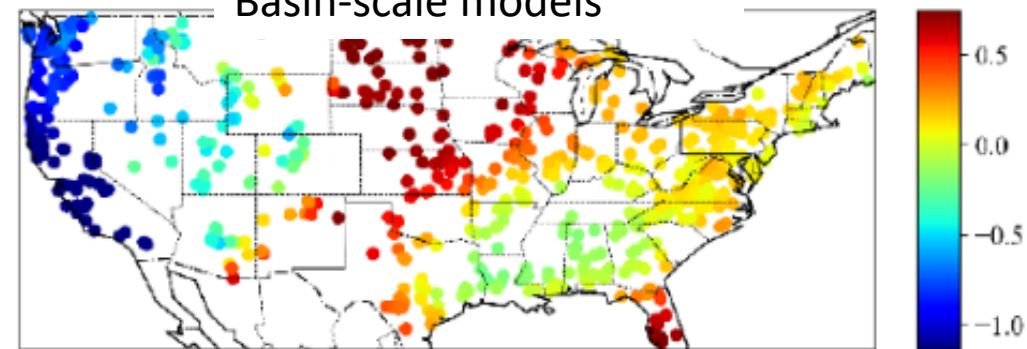
- Soil Moisture Active Passive (SMAP)
  - Launched recently (2015/04)
  - 2~3 days revisit time
  - Senses moisture-dependent top surface soil
- Streamflow modeling
  - Daily data
  - Accompanying attributes
  - With reservoirs, in data-sparse regions
- Dissolved oxygen
- Water temperature
- Sediment
- Snow water equivalent



Gridded models



Basin-scale models



# Long-term projections (first-phase of DL in hydrology)

- Examined comparison with in-situ data & long-term projections

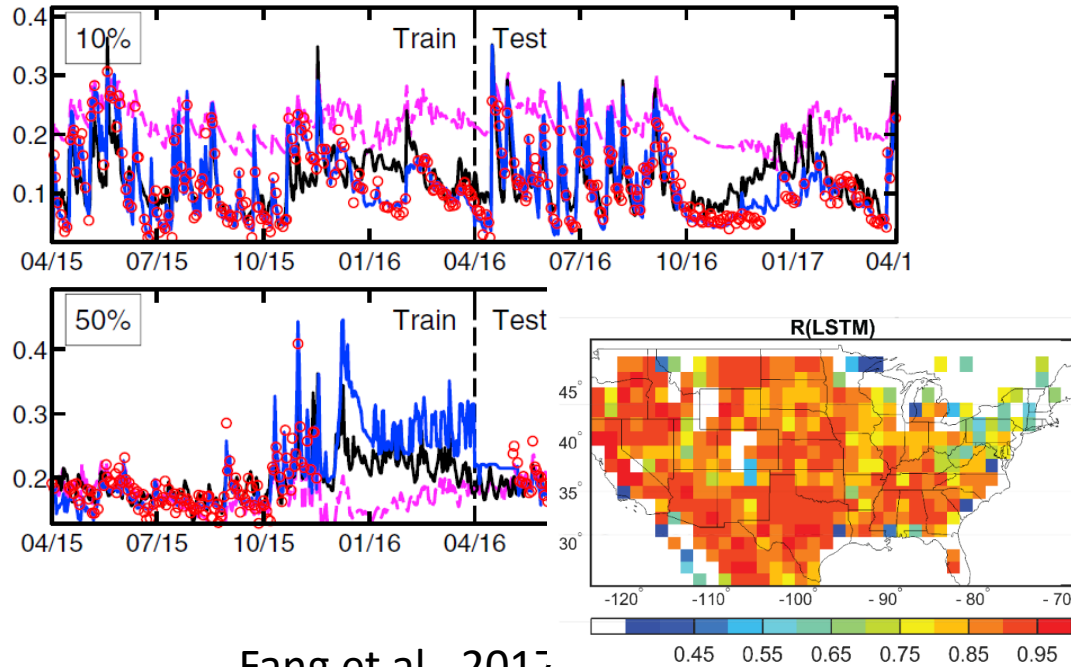
## Geophysical Research Letters

Research Letter | Full Access

Prolongation of SMAP to Spatiotemporally Seamless Coverage of Continental U.S. Using a Deep Learning Neural Network

Kuai Fang, Chaopeng Shen, Daniel Kifer, Xiao Yang

First published: 16 October 2017 | <https://doi.org/10.1002/2017GL075619> | Cited by: 3



Fang et al., 2017

## Water Resources Research

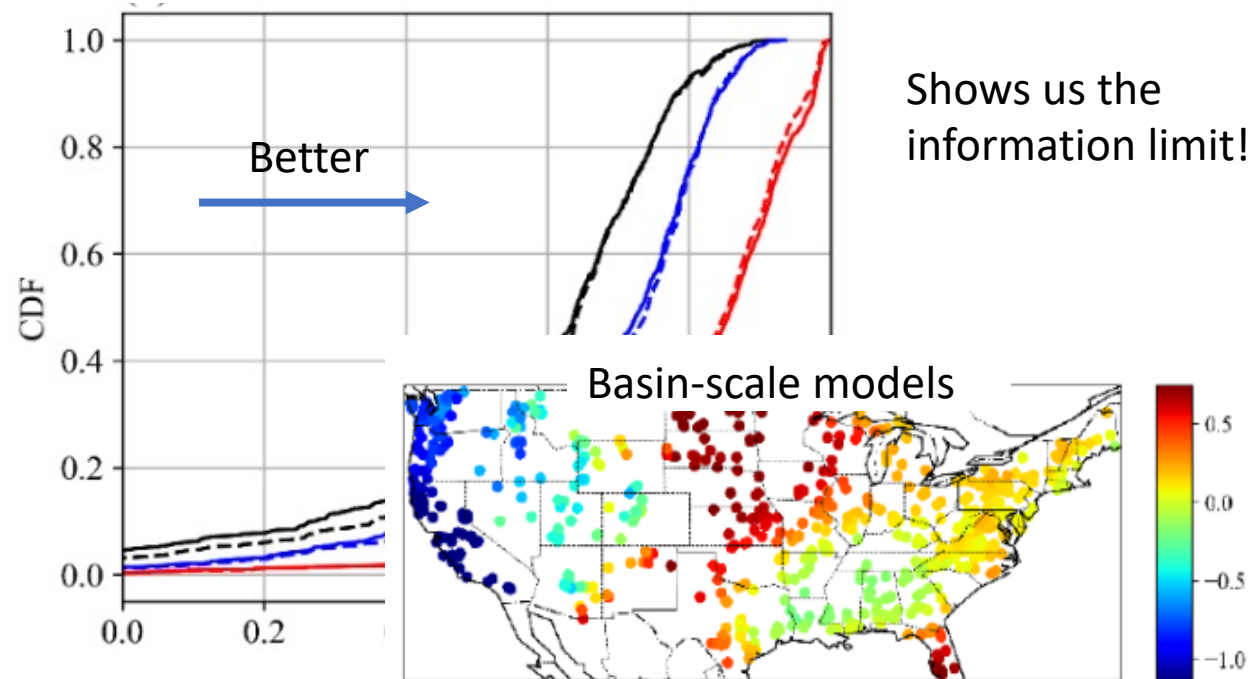
RESEARCH ARTICLE  
10.1029/2019WR026793

Special Section:  
Big Data & Machine Learning  
in Water Sciences: Recent  
Progress and Their Use in  
Advancing Science

Enhancing Streamflow Forecast and Extracting Insights Using Long-Short Term Memory Networks With Data Integration at Continental Scales

Dapeng Feng<sup>1</sup>, Kuai Fang<sup>1,2</sup>, and Chaopeng Shen<sup>1</sup>

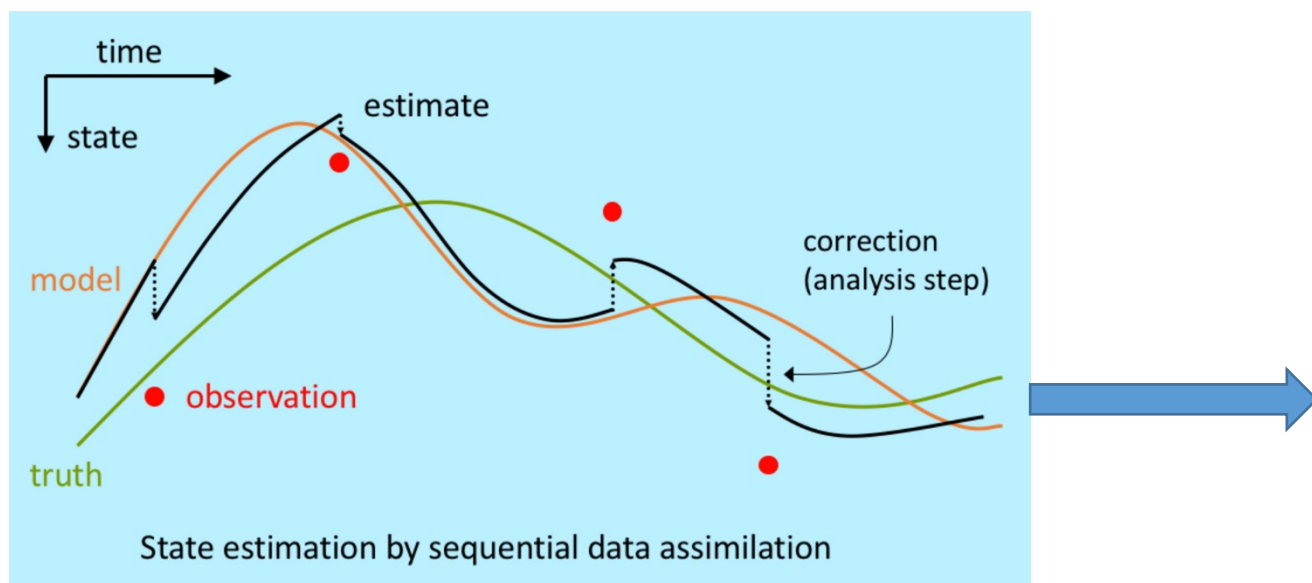
<sup>1</sup>Civil and Environmental Engineering, Pennsylvania State University, State College, PA, USA, <sup>2</sup>Now at: Earth System Science, Stanford University, Stanford, CA, USA





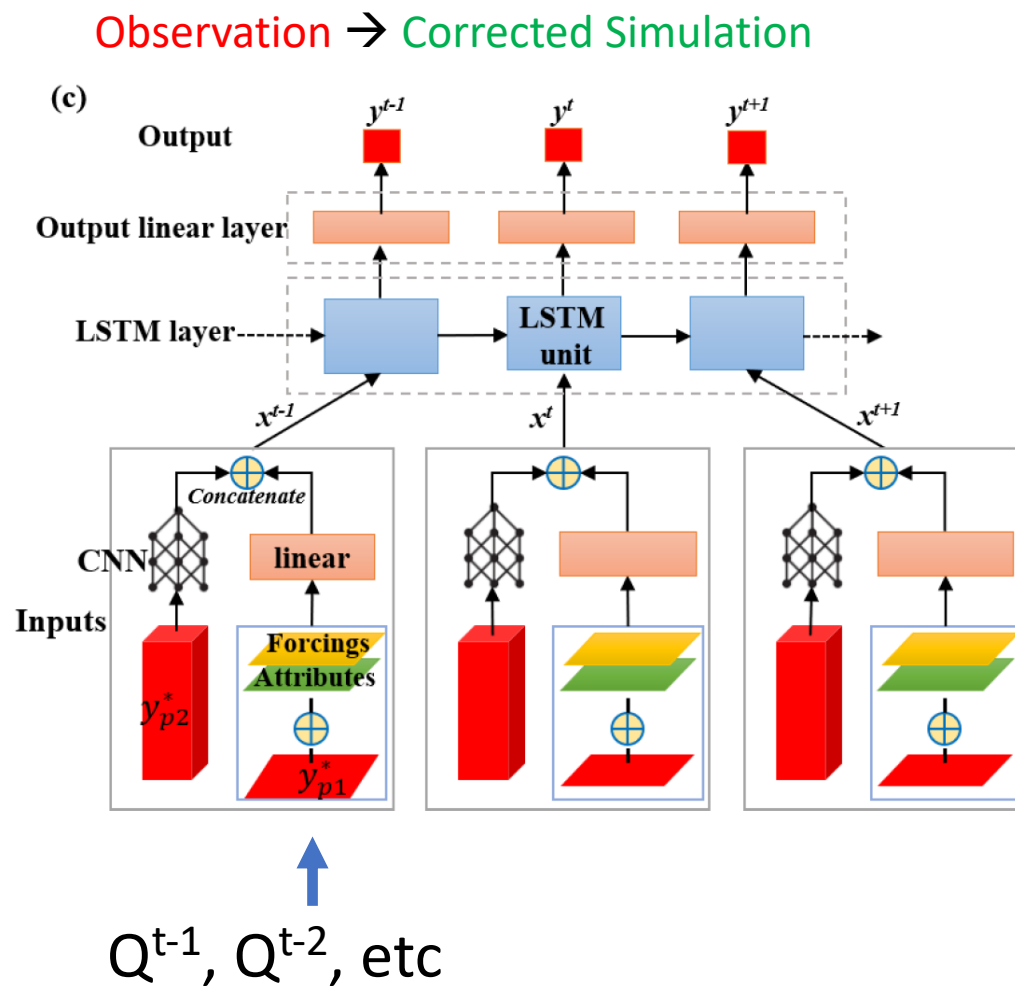
# Short-term forecast

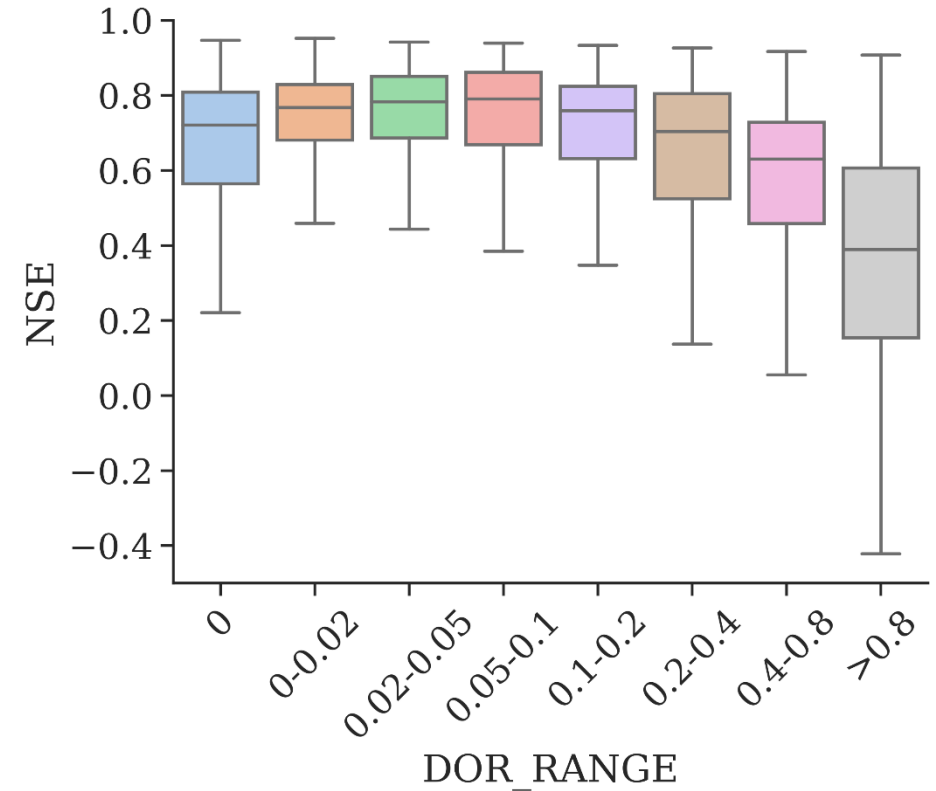
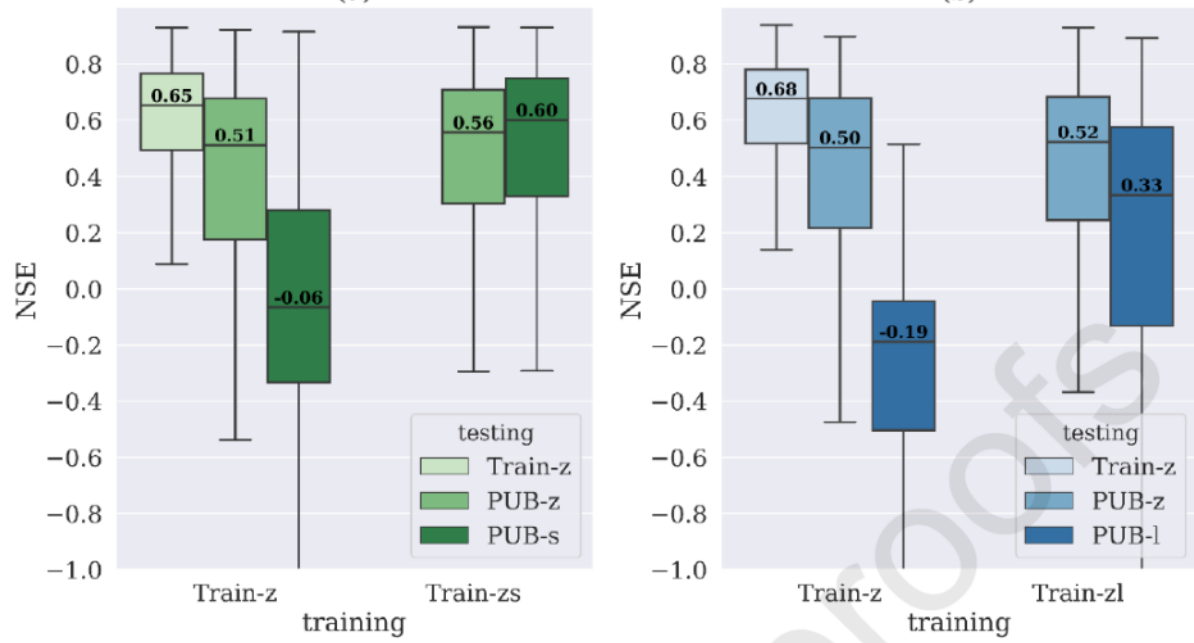
- Traditional “data assimilation” scheme



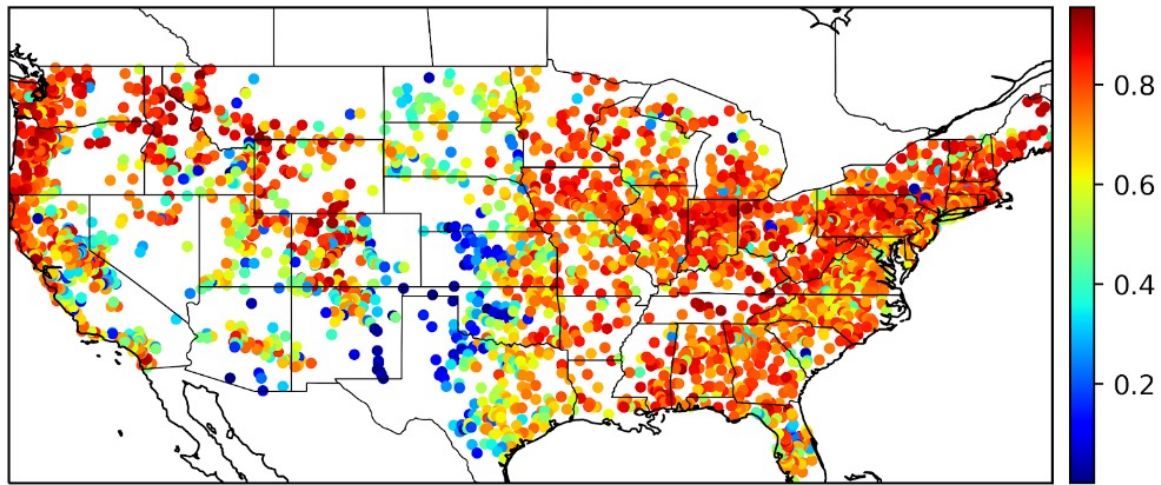
Simulation  $\rightarrow$  Observation – (ENKF)  $\rightarrow$  Correction

Choices: covariance matrix, what to include, how to solve, bias correction, etc.

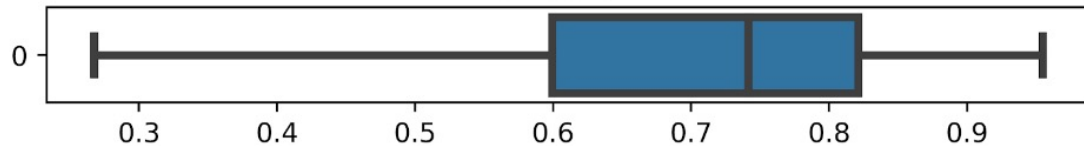




NSE map



NSE boxplot

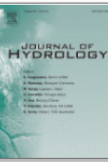


ELSEVIER

Journal of Hydrology

Available online 16 May 2021, 126455

In Press, Journal Pre-proof



Research papers

# Continental-scale streamflow modeling of basins with reservoirs: towards a coherent deep-learning-based strategy

Wenyu Ouyang<sup>a</sup>, Kathryn Lawson<sup>b</sup>, Dapeng Feng<sup>b</sup>, Lei Ye<sup>a</sup>, Chi Zhang<sup>a</sup>, Chaopeng Shen<sup>b</sup>  




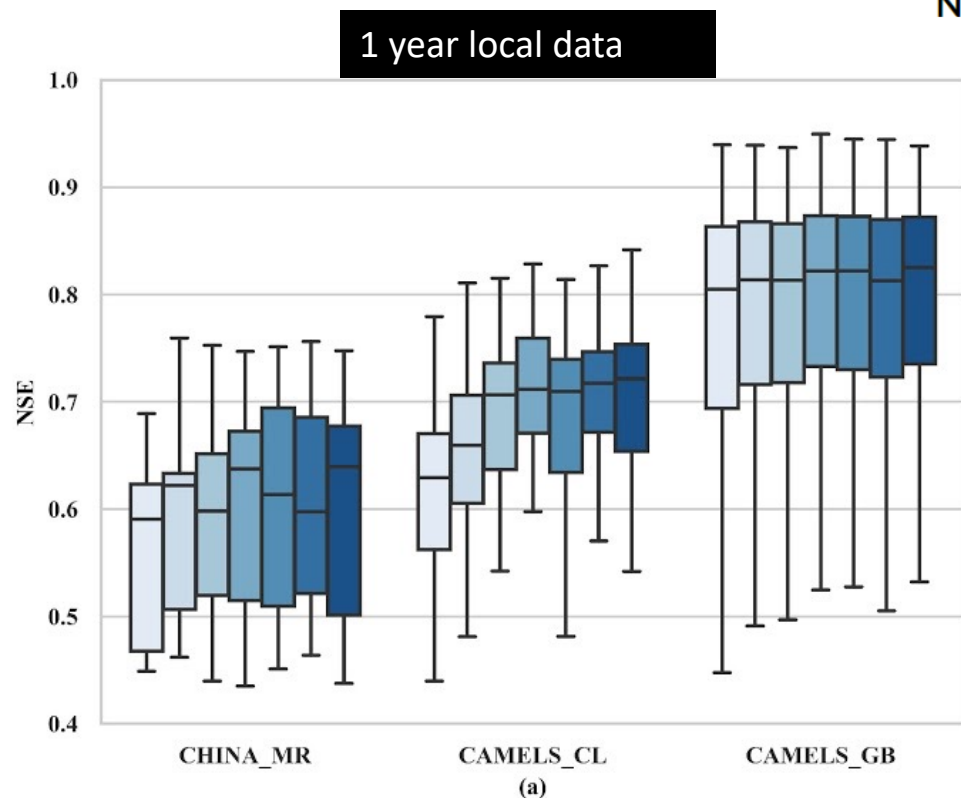
# Sparse-data region

- Transfer learning

## RESEARCH ARTICLE

## How to enhance hydrological predictions in hydrologically distinct watersheds of the Indian subcontinent?

Nikunj K. Mangukiya<sup>1</sup> | Ashutosh Sharma<sup>1</sup>  | Chaopeng Shen<sup>2</sup>



Ma et al., WRR

<https://doi.org/10.1029/2020WR028600>

Basis of

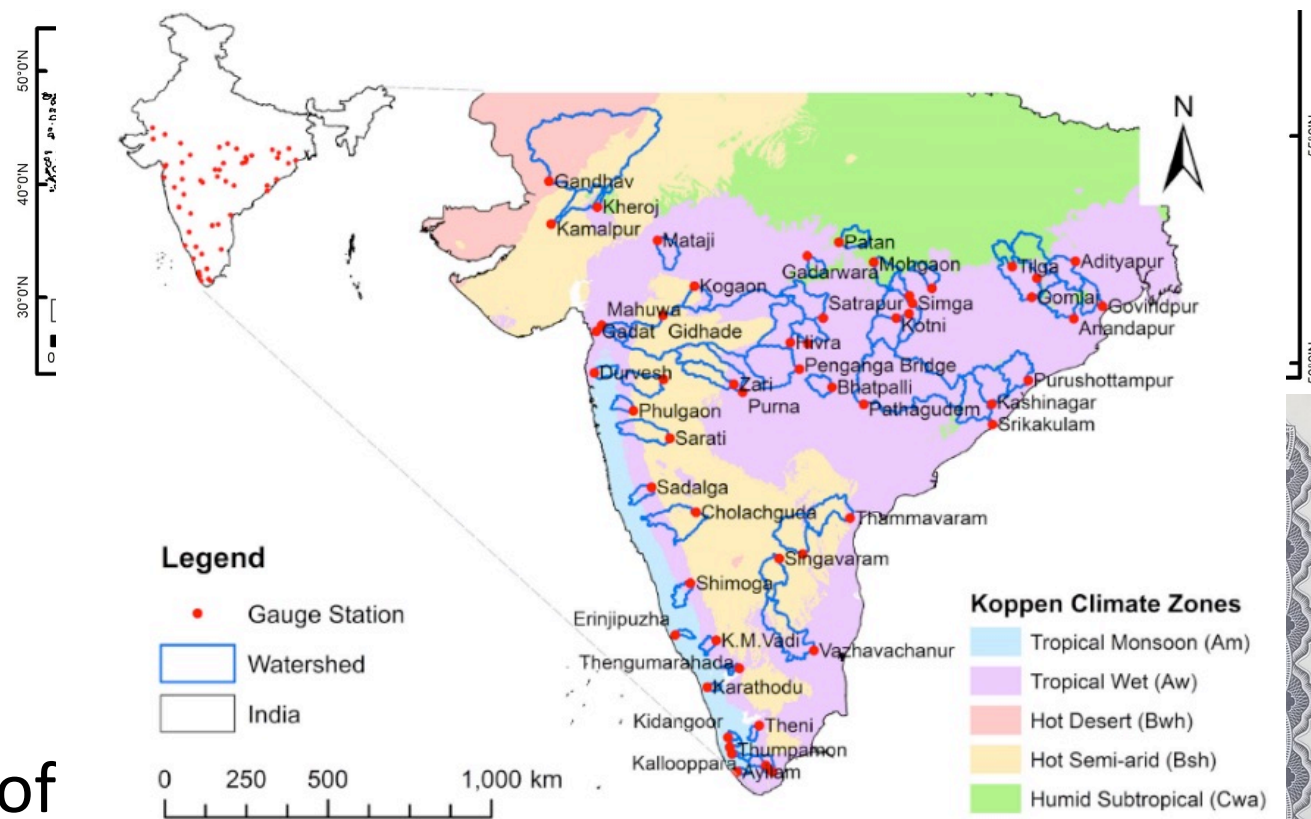
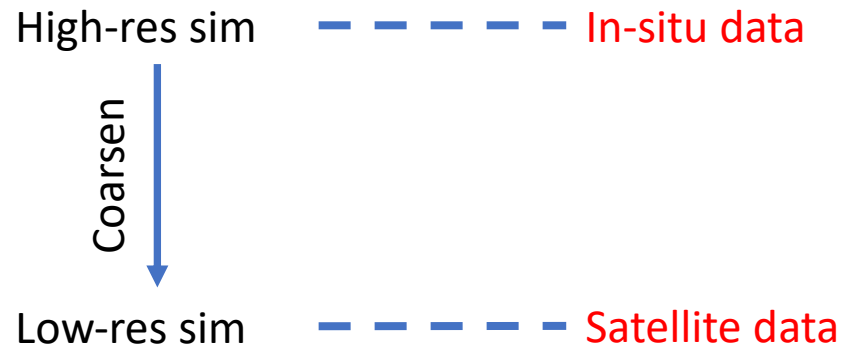
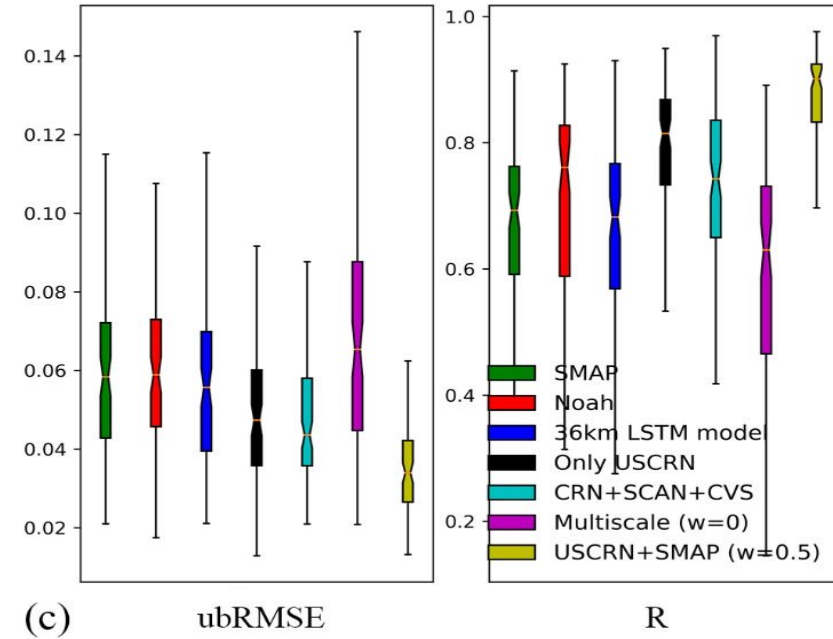
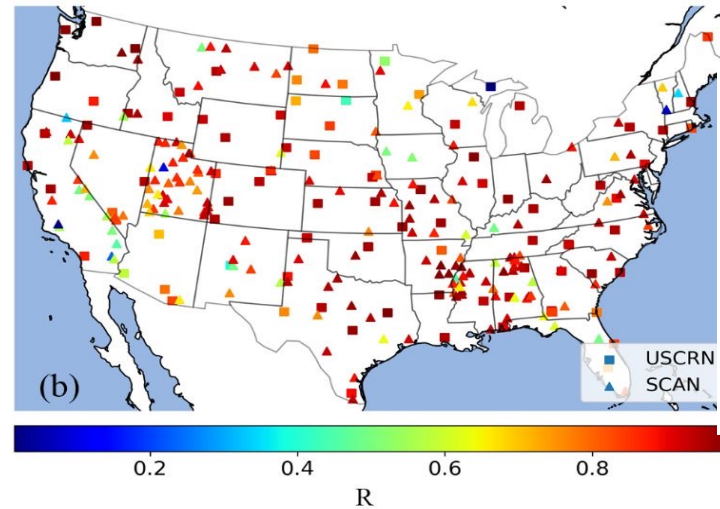
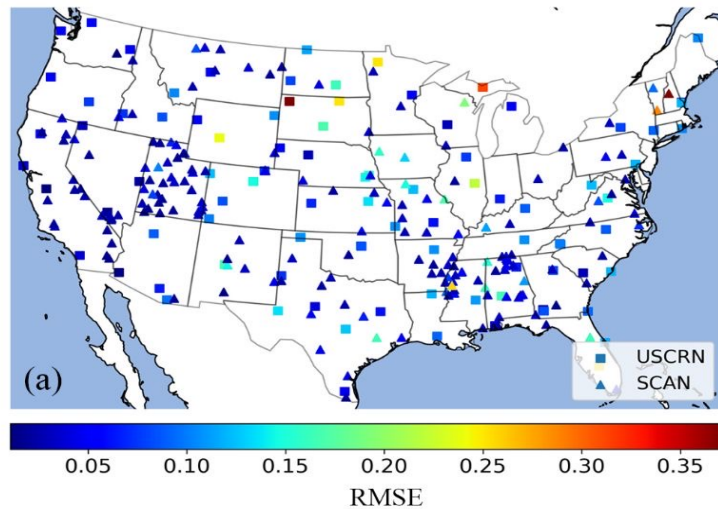


FIGURE 1 Locations of the gauge station and watershed across India. The gauges are spread across different Koppen climate zones and have distinct watershed characteristics.

# Multiscale soil moisture – learning from two teachers



Test period: 2015-04-01 to 2020-03-31



Geophysical Research Letters®

Research Letter | Full Access

A multiscale deep learning model for soil moisture integrating satellite and in-situ data

Jiangtao Liu, Farshid Rahmani, Kathryn Lawson, Chaopeng Shen

First published: 14 March 2022 | <https://doi.org/10.1029/2021GL096847>

# Water quality

## nature water

Explore content ▾ About the journal ▾ Publish with us ▾ Subscribe

[nature](#) > [nature water](#) > [articles](#) > article

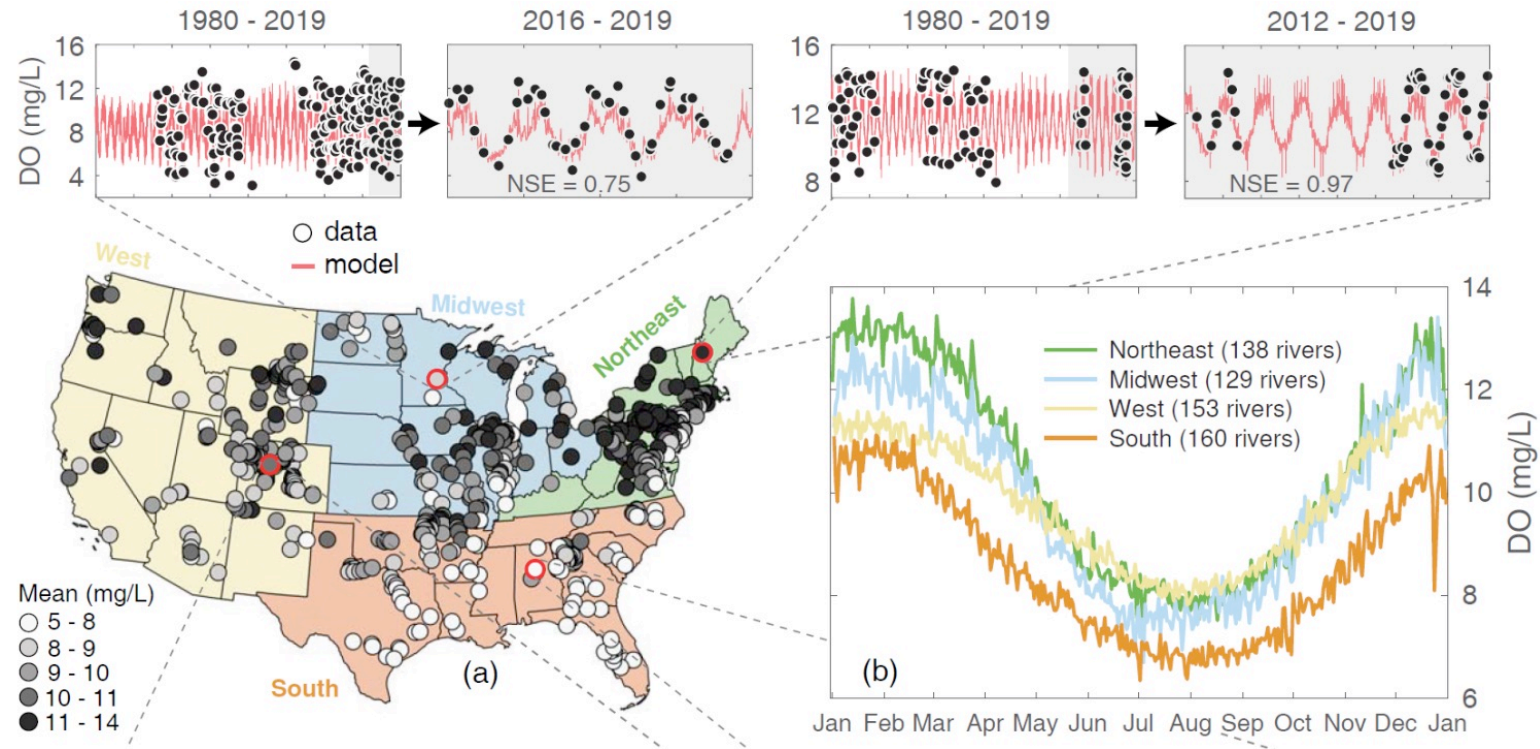
Article | Published: 09 March 2023

### Temperature outweighs light and flow as the predominant driver of dissolved oxygen in US rivers

[Wei Zhi](#), [Wenyu Ouyang](#), [Chaopeng Shen](#) & [Li Li](#)

[Nature Water](#) 1, 249–260 (2023) | [Cite this article](#)

## Dissolved Oxygen



ELSEVIER

Contents lists available at [ScienceDirect](#)

Science of the Total Environment

journal homepage: [www.elsevier.com/locate/scitotenv](http://www.elsevier.com/locate/scitotenv)



A deep learning-based novel approach to generate continuous daily stream nitrate concentration for nitrate data-sparse watersheds

Gourab Kumer Saha<sup>a</sup>, Farshid Rahmani<sup>b</sup>, Chaopeng Shen<sup>b</sup>, Li Li<sup>b</sup>, Raj Cibin<sup>a,b,\*</sup>

<sup>a</sup> Department of Agricultural and Biological Engineering, The Pennsylvania State University, United States of America

<sup>b</sup> Department of Civil and Environmental Engineering, The Pennsylvania State University, United States of America

## Nitrate

## ENVIRONMENTAL RESEARCH LETTERS

### LETTER

### Exploring the exceptional performance of a deep learning stream temperature model and the value of streamflow data

Farshid Rahmani<sup>1</sup>, Kathryn Lawson<sup>1</sup>, Wenyu Ouyang<sup>2</sup>, Alison Appling<sup>3</sup>, Samantha Oliver<sup>4</sup> and Chaopeng Shen<sup>1</sup>

<sup>1</sup> Civil and Environmental Engineering, Pennsylvania State University, University Park, State College, PA, United States of America

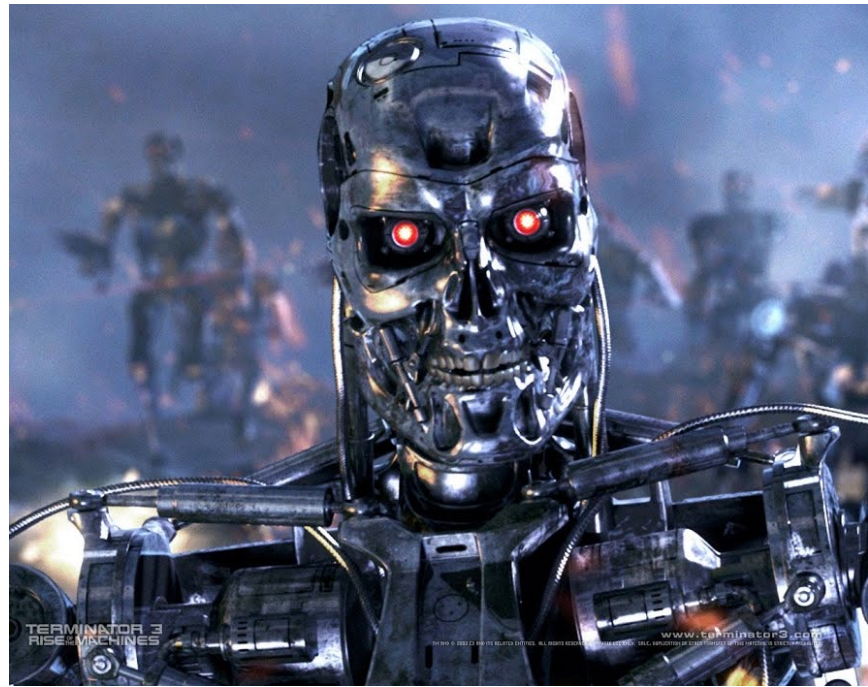
<sup>2</sup> School of Hydraulic Engineering, Dalian University of Technology, Dalian, People's Republic of China

<sup>3</sup> US Geological Survey, Reston, VA, United States of America

<sup>4</sup> US Geological Survey, Upper Midwest Water Science Center, Middleton, WI, United States of America

## Water temperature



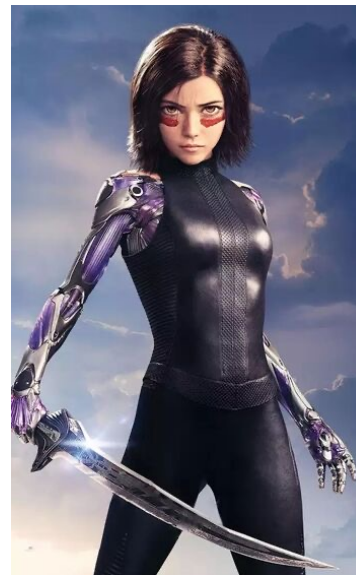


pure DL models



Human modelers

We need hybrids



# Phase 2: How to surpass the teacher (training data)

Training data often have limitations:

Resolution, accuracy, time interval, availability (unobserved variables), geographical imbalance, not enough extremes, not capturing nonstationarity...

How to overcome such limitations?

- Inclusion of physics
- Learning about physics.



# *Similarity & Differences* between deep learning (DL) and process-based models (PBM)?



[This Photo](#) by Unknown Author is licensed under [CC BY-NC](#)

Purely data-driven NNs	Purely process-based models
<b>Similarities</b>	
$y = g^W(u, x, A)$ $W = \operatorname{argmin}(L(y, y^*))$	$y = f^\theta(u, x, A)$ $\theta = \operatorname{argmin}(L(y, y^*))$



[This Photo](#) by Unknown Author is licensed under [CC BY-SA](#)

The secret? Differentiable programming!



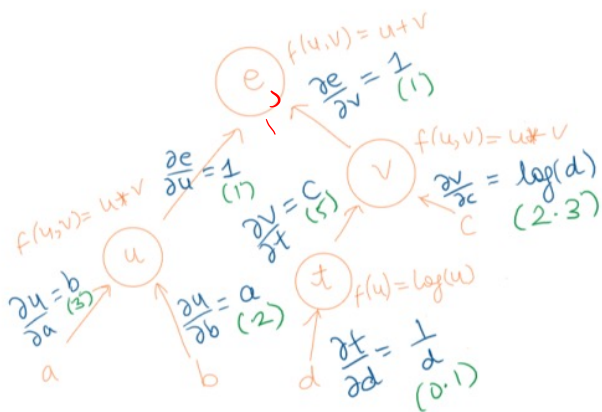
# What does “Differentiable” mean?

- The ability to rapidly compute gradients  $\frac{dL}{d\theta}$
- Enabling training by gradient descent

## Automatic differentiation

Back-propagation:

e.g.  $a=2, b=3, c=5, d=10$



$$e = a * b + c \log(d)$$

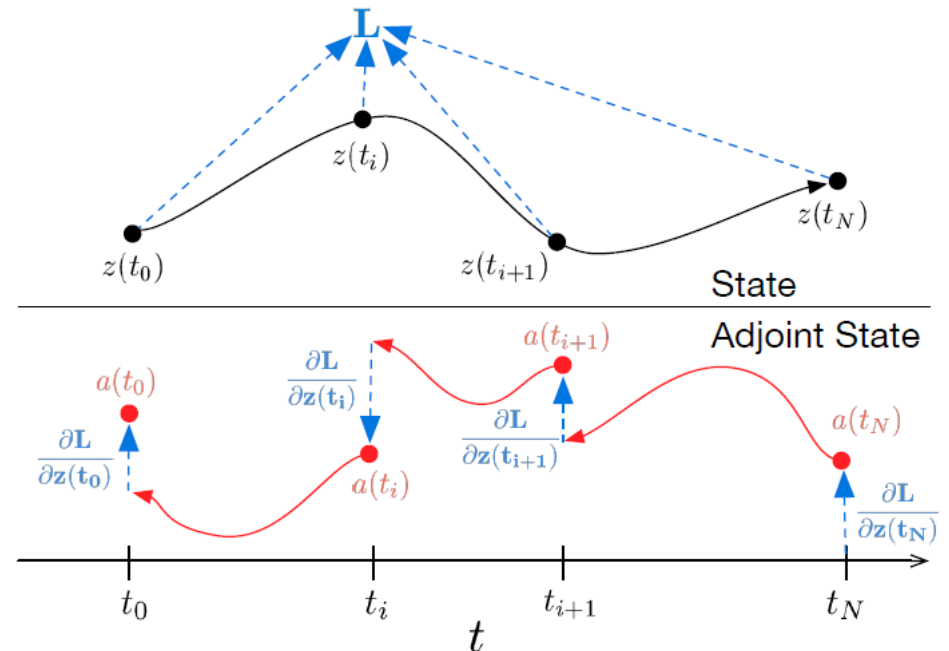
$$\frac{\partial e}{\partial a} = b(1) = b = 3$$

$$\frac{\partial e}{\partial b} = a(1) = a = 2$$

$$\frac{\partial e}{\partial c} = \log d \times 1 = \log d = 2.3$$

$$\frac{\partial e}{\partial d} = \frac{1}{d} \times c \times 1 = \frac{c}{d} = 0.5$$

## Adjoint State method



# Differentiable parameter learning



ARTICLE

<https://doi.org/10.1038/s41467-021-26107-z>

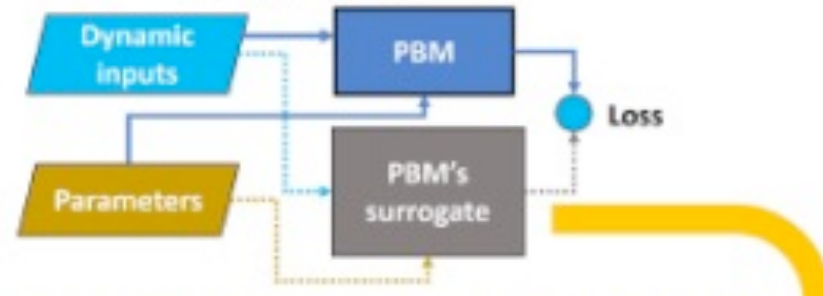
OPEN



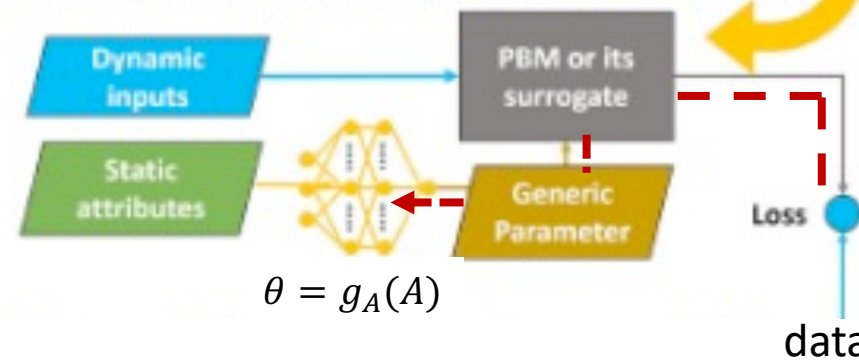
From calibration to parameter learning: Harnessing the scaling effects of big data in geoscientific modeling

Wen-Ping Tsai<sup>1</sup>, Dapeng Feng<sup>1</sup>, Ming Pan<sup>2,3</sup>, Hylke Beck<sup>4</sup>, Kathryn Lawson<sup>1,5</sup>, Yuan Yang<sup>6,7</sup>, Jiangtao Liu<sup>1</sup> & Chaopeng Shen<sup>1,5</sup>

(a) PBM or PBM's surrogate (optional)

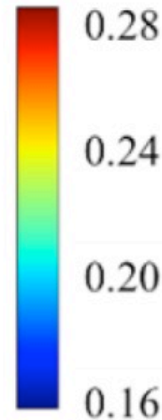
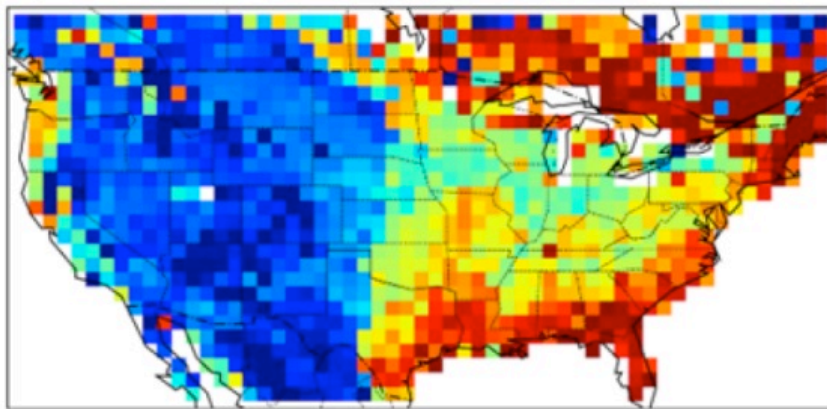


(b) dPL  $g_A$  framework (if historical observations are unavailable)

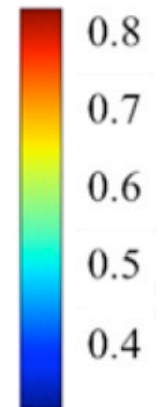
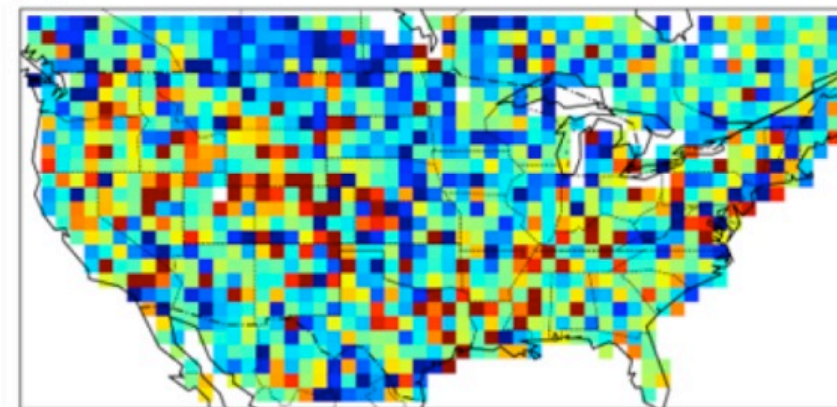


$$\theta = g_A(A)$$

(a) dPL  $g_z$  INFILT



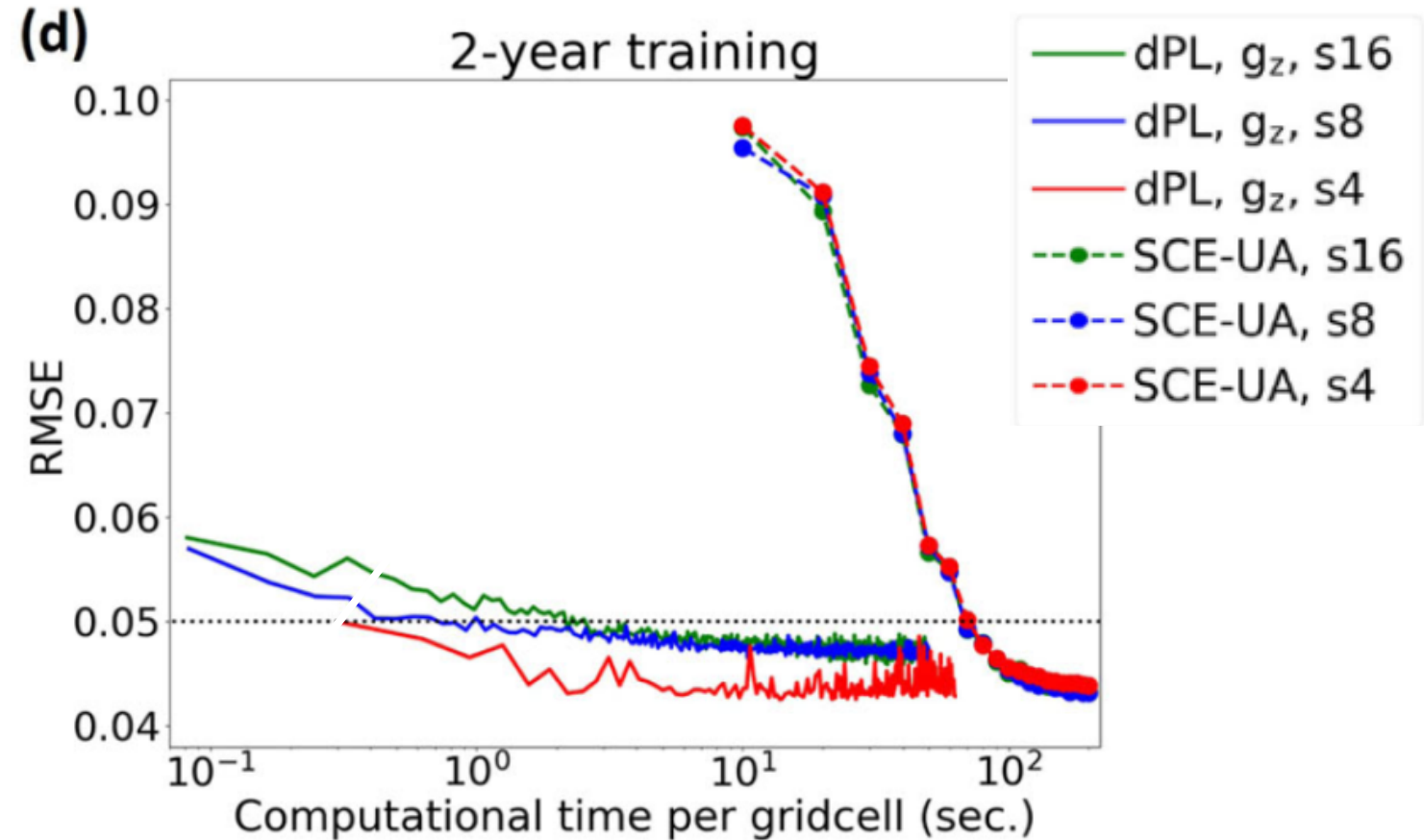
(b) SCE INFILT



# Point #1. Data scaling relationships (network effect?)

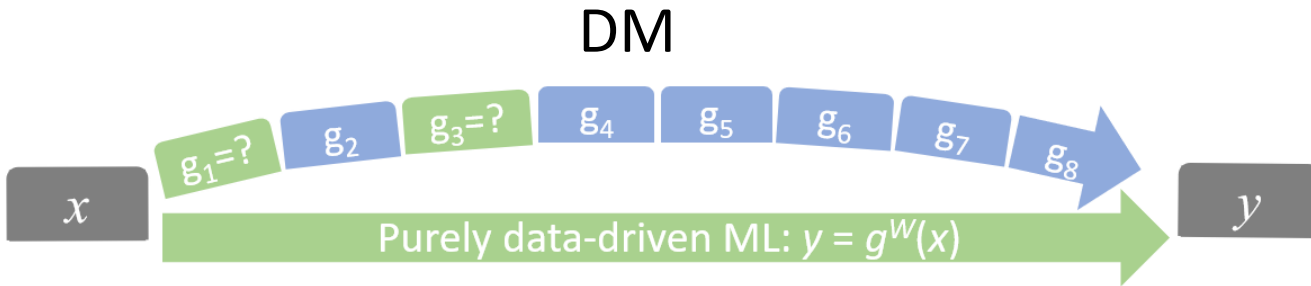
1. dPL = SCEUA for lowest RMSE
2. dPL scales better with more data
3. Orders of magnitude more efficient
4. (not shown) better results for **untrained** variables and better **spatial generalization** than traditional approach!

Relies on differentiable programming!

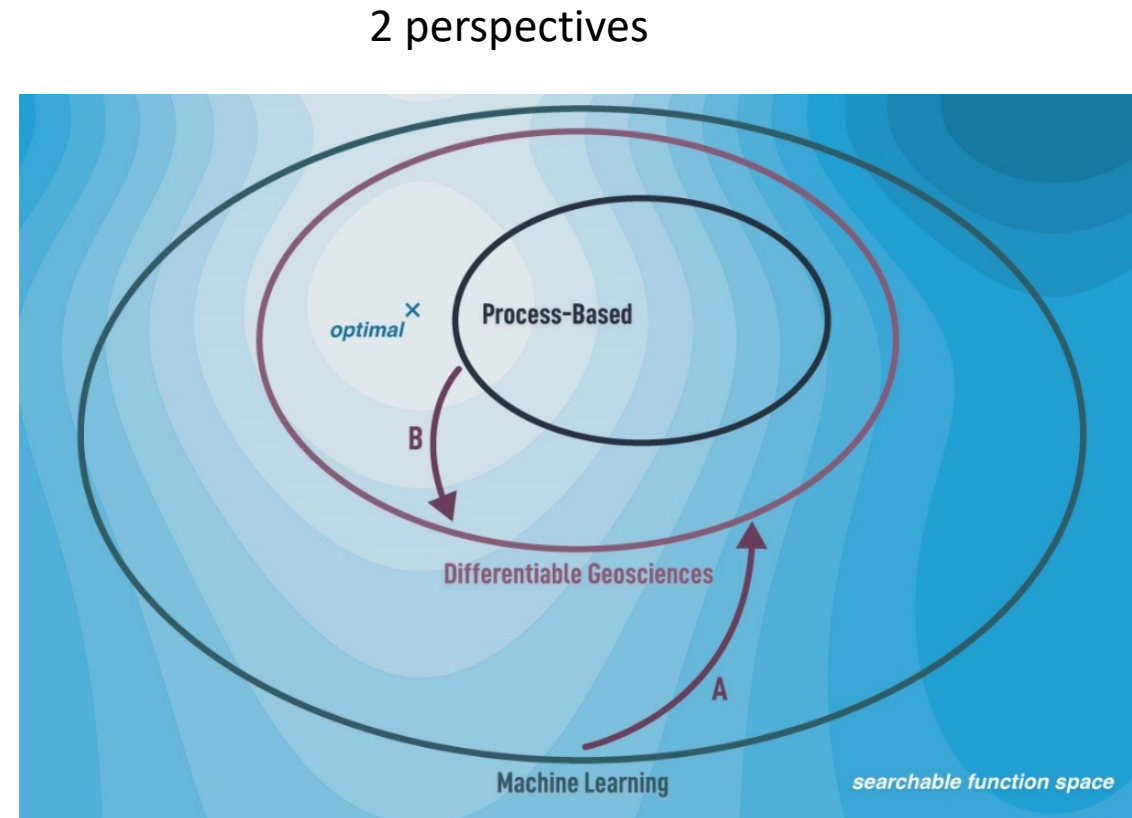




# What is Differentiable Modeling (DM) in Geosciences?



- NNs mixed w/ process-based equations (priors)
- The priors constrain the learning to an interpretable scope.
- intermediate physical variables.
- Update our knowledge and learn unrecognized relationships from data.



# Differentiable, learnable models to learn functions

Hydrol. Earth Syst. Sci., 27, 2357–2373, 2023  
<https://doi.org/10.5194/hess-27-2357-2023>  
 © Author(s) 2023. This work is distributed under the Creative Commons Attribution 4.0 License.



The suitability of differentiable, physics-informed machine learning hydrologic models for ungauged regions and climate change impact assessment

Dapeng Feng<sup>1</sup>, Hylke Beck<sup>2</sup>, Kathryn Lawson<sup>1</sup>, and Chaopeng Shen<sup>1</sup>

<sup>1</sup>Civil and Environmental Engineering, The Pennsylvania State University, University Park, PA, USA

<sup>2</sup>Physical Science and Engineering, King Abdullah University of Science and Technology, Thuwal, Saudi Arabia

Correspondence: Chaopeng Shen (cshen@engr.psu.edu)

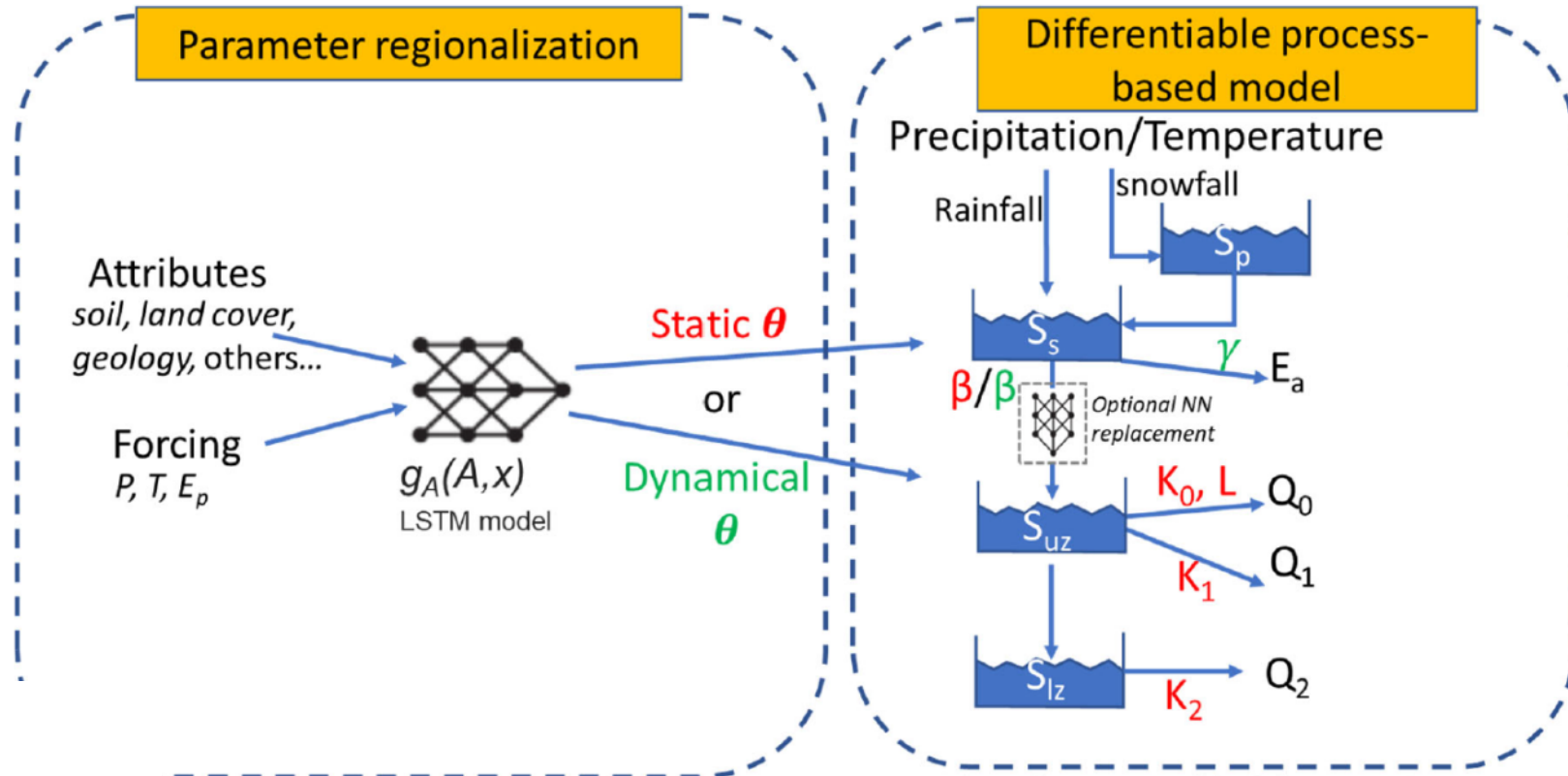
## Water Resources Research

Research Article | [Full Access](#)

Differentiable, learnable, regionalized process-based models with multiphysical outputs can approach state-of-the-art hydrologic prediction accuracy

Dapeng Feng, Jiangtao Liu, Kathryn Lawson, Chaopeng Shen

First published: 19 September 2022 | <https://doi.org/10.1029/2022WR032404>



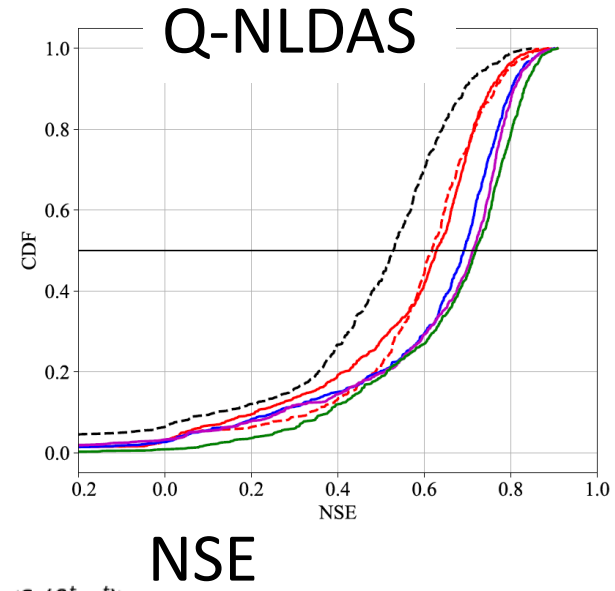
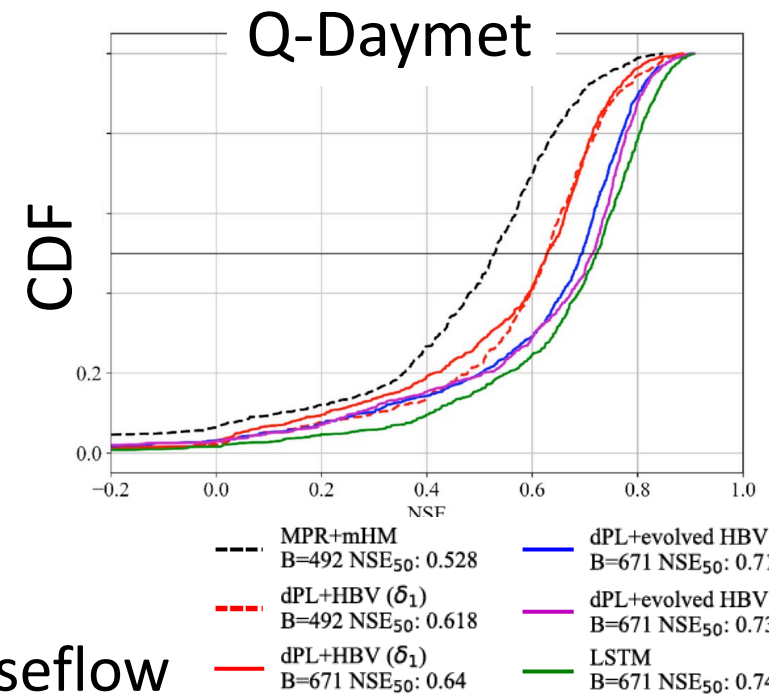
\* Not all parameters and detailed processes of HBV sketched here for the sake of simplicity.

Rewritten in PyTorch

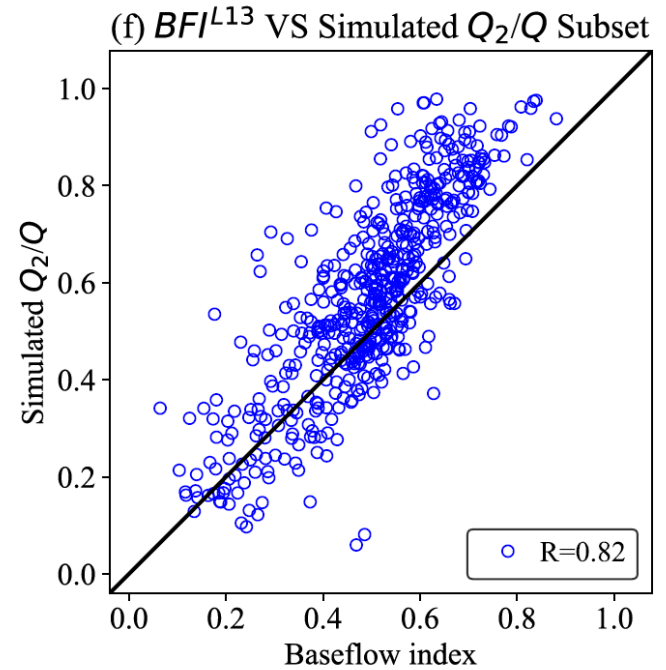
*Evolve model structure*

Approaching LSTM!  
But....

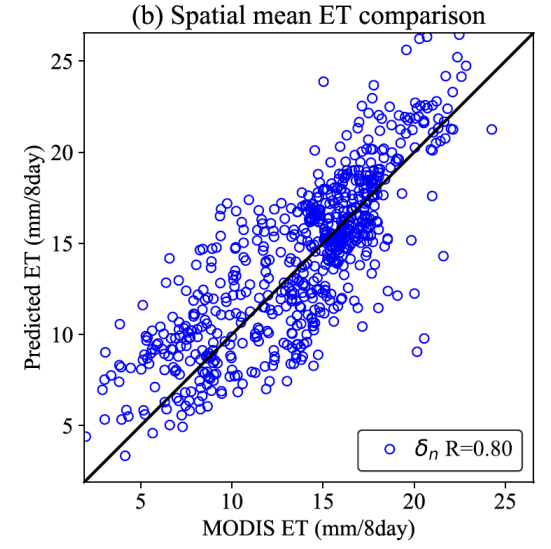
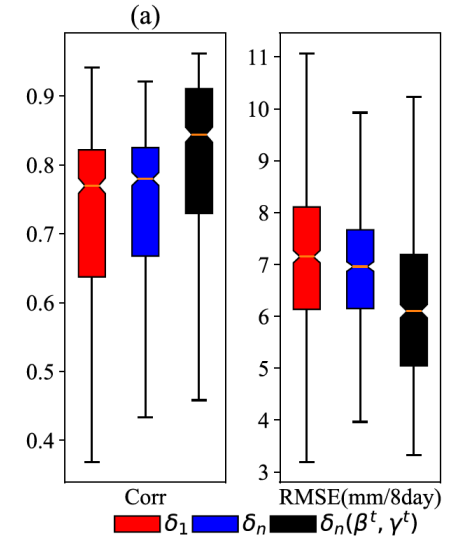
- Output untrained variables.
- Multivariate constraints.
- It can help us answer questions!



Baseflow

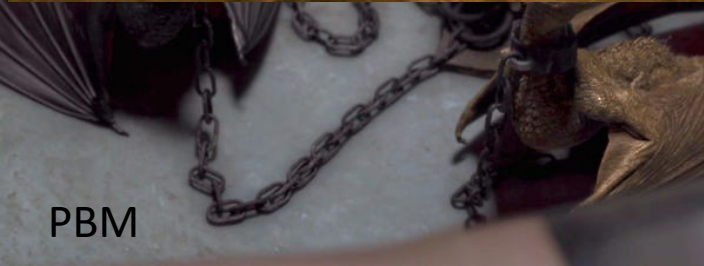


Evapotranspiration



Caveat: not using the ensemble  
-- first iteration. Priors do matter.





PBM

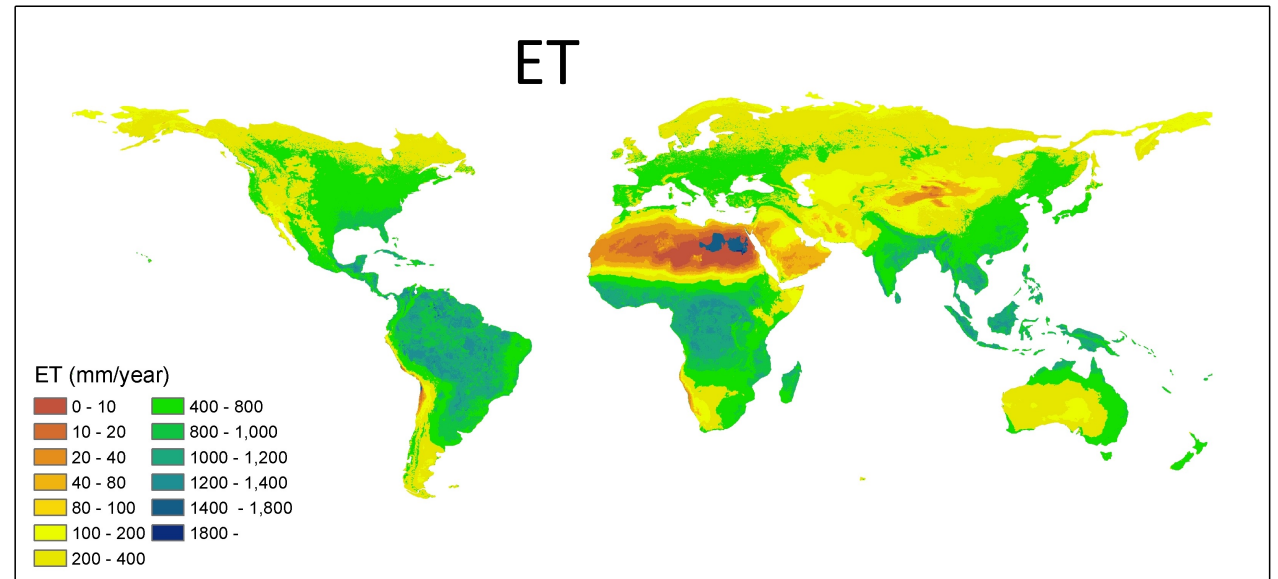
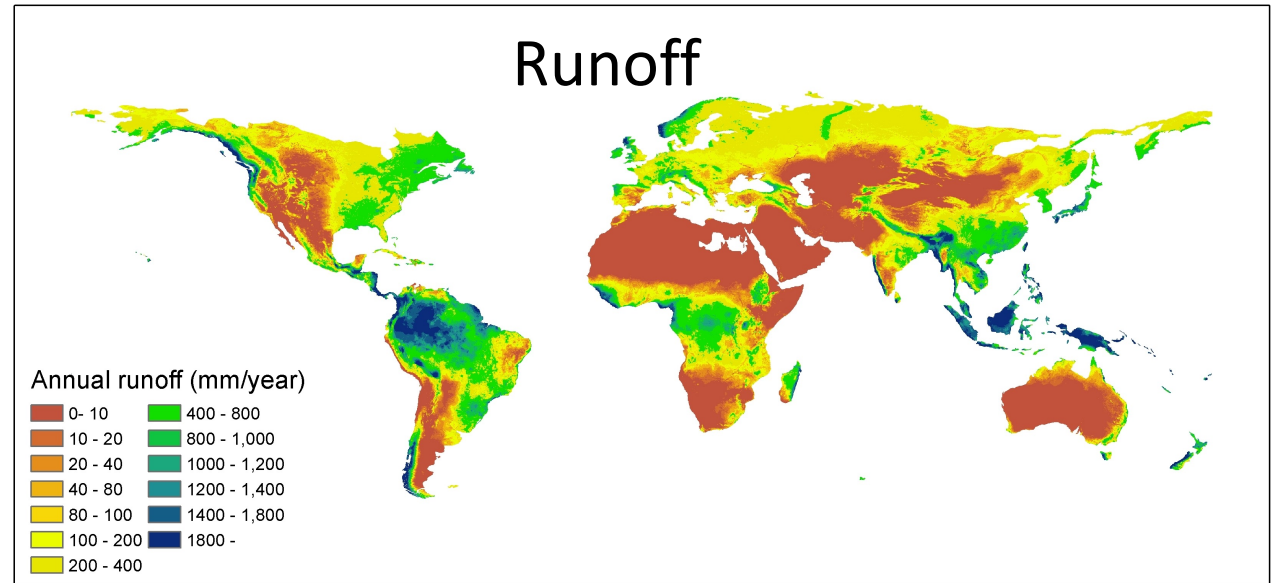


PBM+dPL

# What can DM bring to global hydrology?

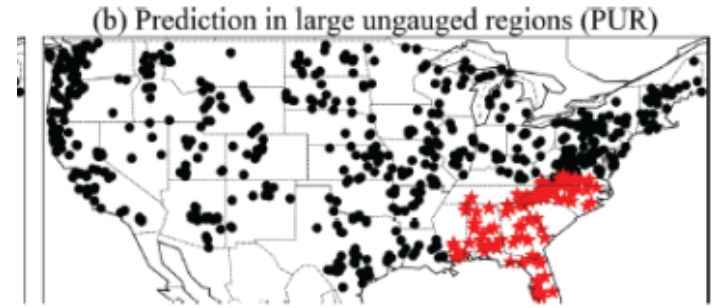
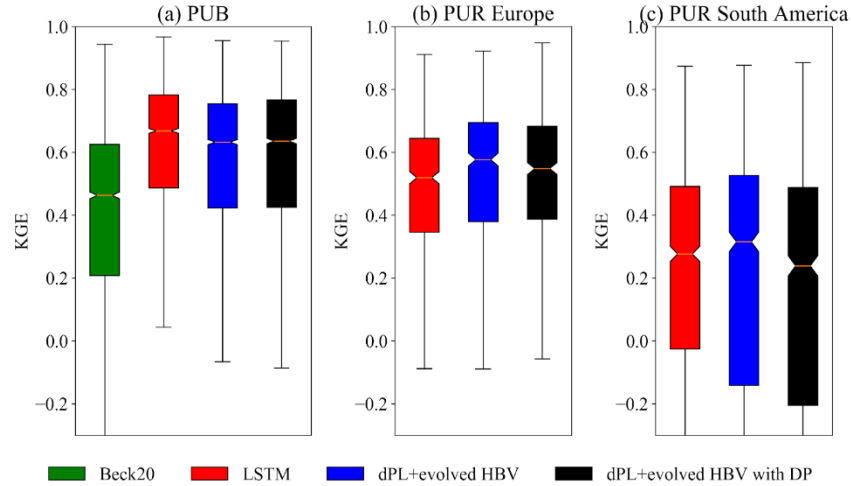
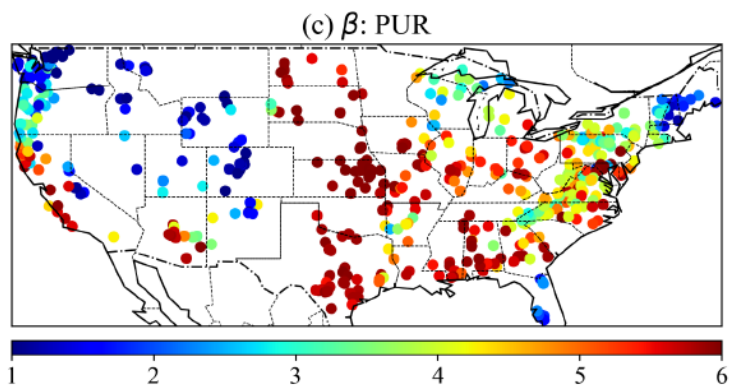
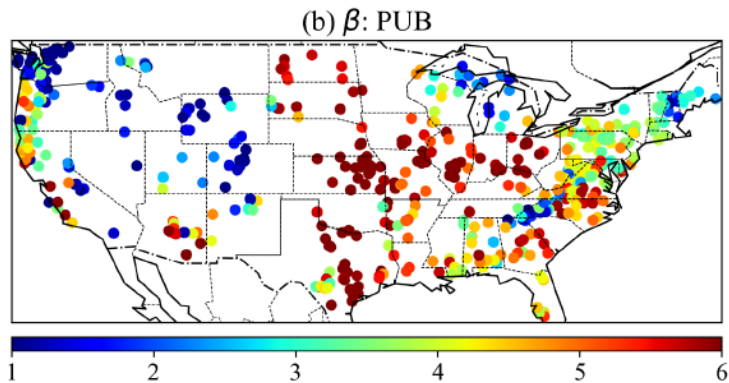
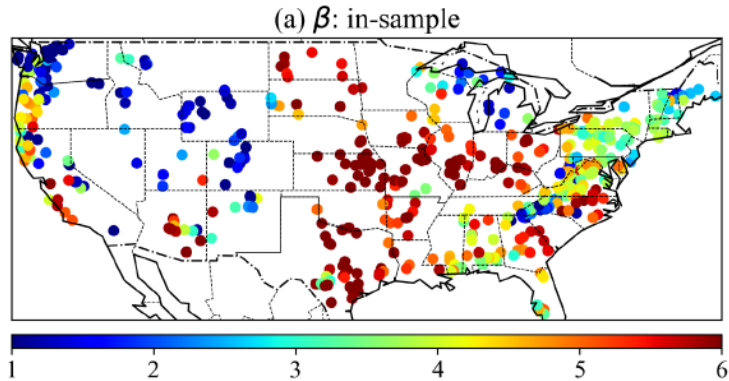
- Spatial extrapolation in data-sparse regions
- Extremes
- Learn robust unknown functions
- Human dynamics or unknown physics
- Correct forcings

Produced by differentiable models

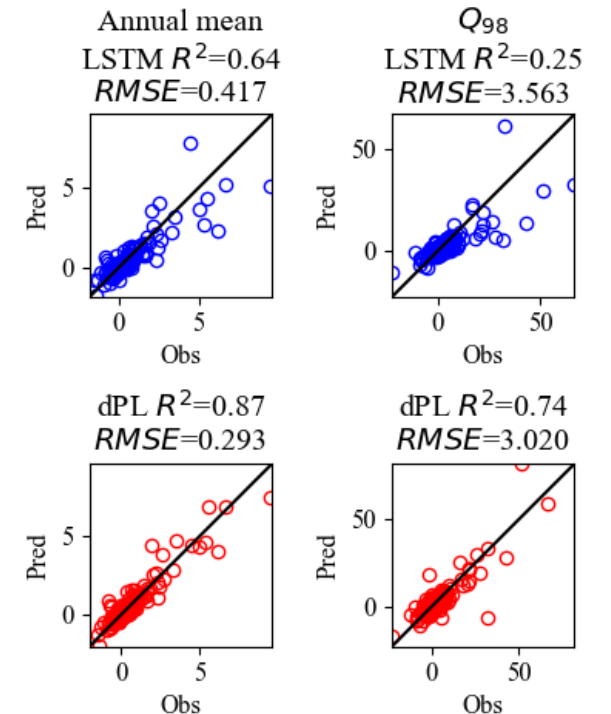
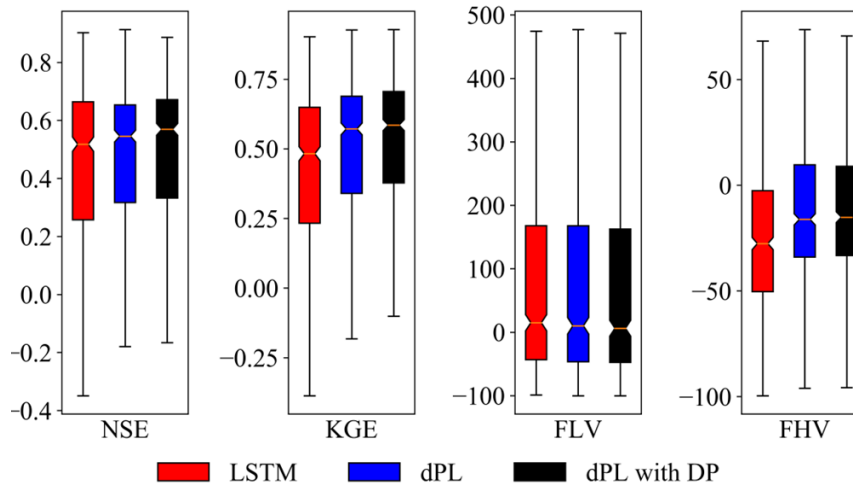




# Differentiable models extrapolation better



## For PUR





## Differentiable, Learnable, Regionalized Process-Based Models With Multiphysical Outputs can Approach State-Of-The-Art Hydrologic Prediction Accuracy

Dapeng Feng, Jiangtao Liu, Kathryn Lawson, Chaopeng Shen

First published: 19 September 2022 | <https://doi.org/10.1029/2022WR032404> | Citations: 18

<https://doi.org/10.5194/hess-27-2357-2023>  
© Author(s) 2023. This work is distributed under the Creative Commons Attribution 4.0 License.

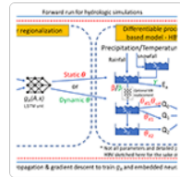
Article Assets Peer review Metrics Related articles

Research article |

30 Jun 2023

The suitability of differentiable, physics-informed machine learning hydrologic models for ungauged regions and climate change impact assessment

Dapeng Feng, Hylke Beck, Kathryn Lawson, and Chaopeng Shen



<https://doi.org/10.5194/gmd-2023-190>  
© Author(s) 2023. This work is distributed under the Creative Commons Attribution 4.0 License.

Abstract Discussion Metrics

Submitted as: model evaluation paper |

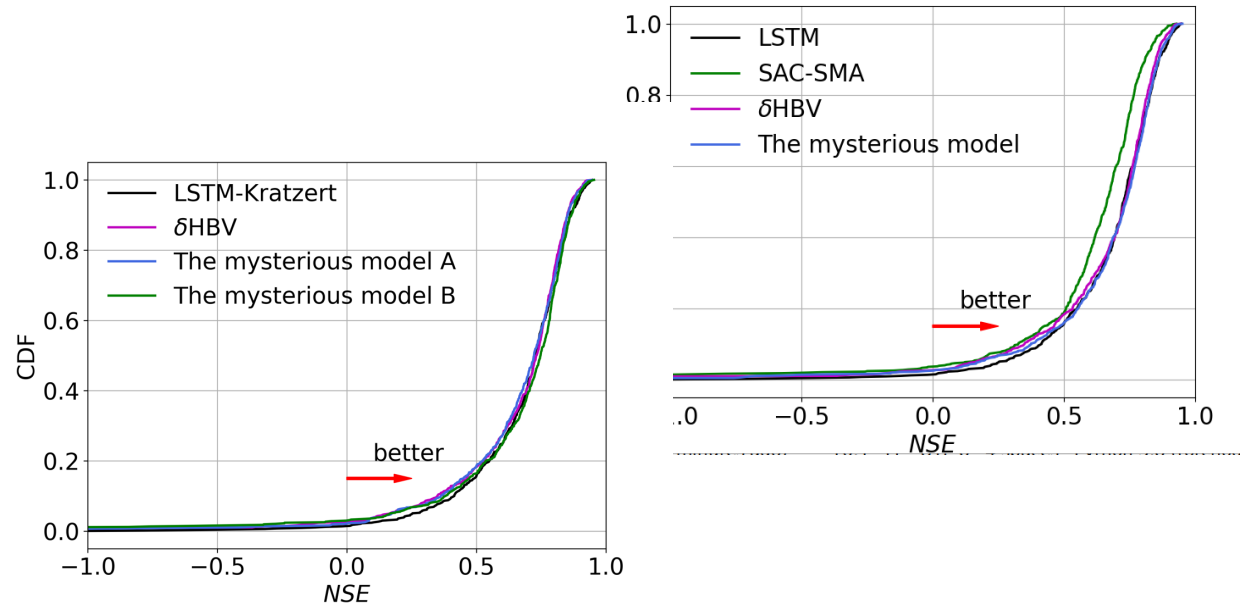
05 Oct 2023

Status: this preprint is currently under review for the journal GMD.

Deep Dive into Global Hydrologic Simulations: Harnessing the Power of Deep Learning and Physics-informed Differentiable Models ( $\delta$ HBV-globe1.0-hydroDL)

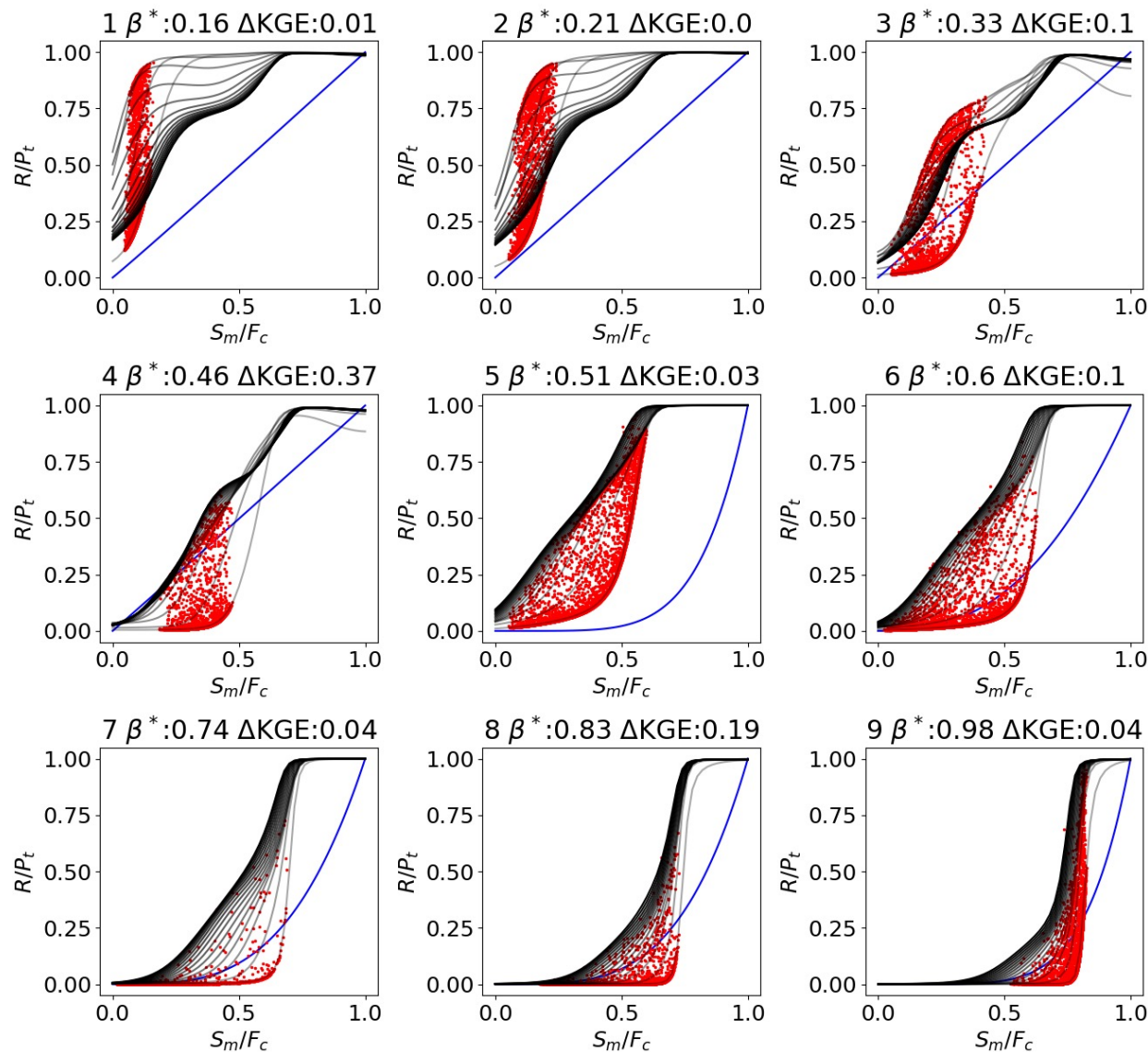
Dapeng Feng, Hylke Beck, Jens de Bruijn, Reetik Kumar Sahu, Yusuke Satoh, Yoshihide Wada, Jiangtao Liu, Ming Pan, Kathryn Lawson, and Chaopeng Shen

# New model



Model	Median NSE	Median KGE	Median absolute (non-absolute) FLV (%)	Median absolute (non-absolute) FHV (%)	Median low flow RMSE (mm/day)	Median peak flow RMSE (mm/day)	Baseflow index spatial correlation	Median NSE of temporal ET simulation
LSTM	<b>0.73</b>	<b>0.77</b>	40.59 (29.70)	13.46 (-4.19)	0.055	2.56	-	-
SAC-SMA	0.66	0.73	59.40 (46.96)	17.55 (-9.79)	0.081	3.19	-	-
HBV	<b>0.73</b>	0.73	56.53 (50.93)	15.29 (-8.89)	0.074	2.56	0.76	0.59
The mysterious model	0.72	0.75	43.29 (37.61)	<b>13.25 (-4.33)</b>	<b>0.048</b>	<b>2.47</b>	<b>0.83</b>	<b>0.61</b>

# Learning unknown relationships from data (in preparation)



$$R/P_t = (S_m/F_c)^\beta$$

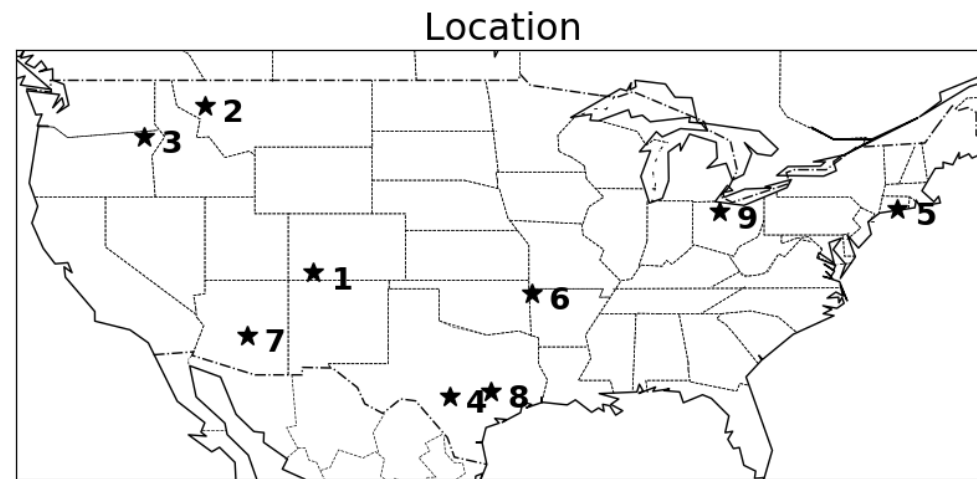


$$R/P_t = ANN(\beta^*, F_c, S_m, S_m/F_c, P_t)$$

Blue line: original power law relation

Red dots: ANN simulations

Black lines: continuous plotting of ANN functions

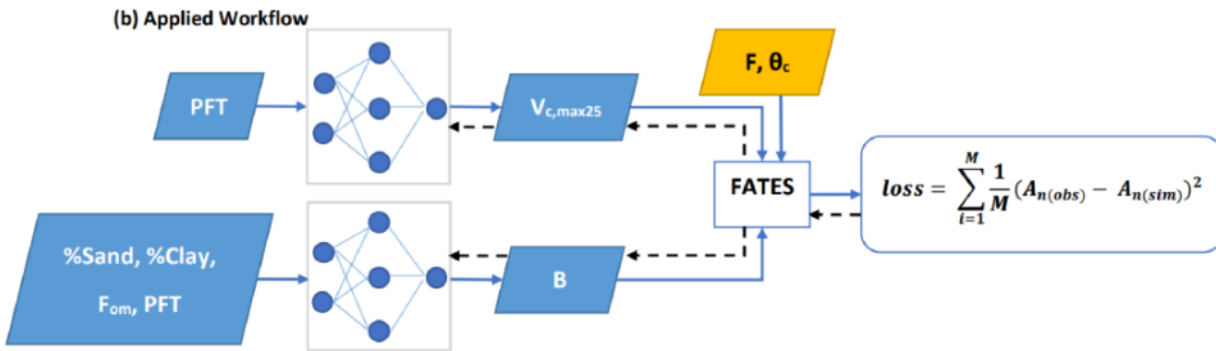


# How is differentiable modeling different from physics-guided ML?

	Physics-guided ML	Differentiable modeling
<b>Goal</b>	Use physics to constrain ML → Improve ML generalization	Use ML to learn unknown relationships and improve simulation quality → Advance our process understanding
<b>Approach</b>	May not be differentiable; Various approaches like modifying the loss function	Differentiable; end-to-end training
<b>Philosophy</b>	Physical laws is treated as ground truth	Constantly seeking to improve our equations



# Example 4. Ecosystem modeling: photosynthesis



(a) Temporal holdout test for the following system

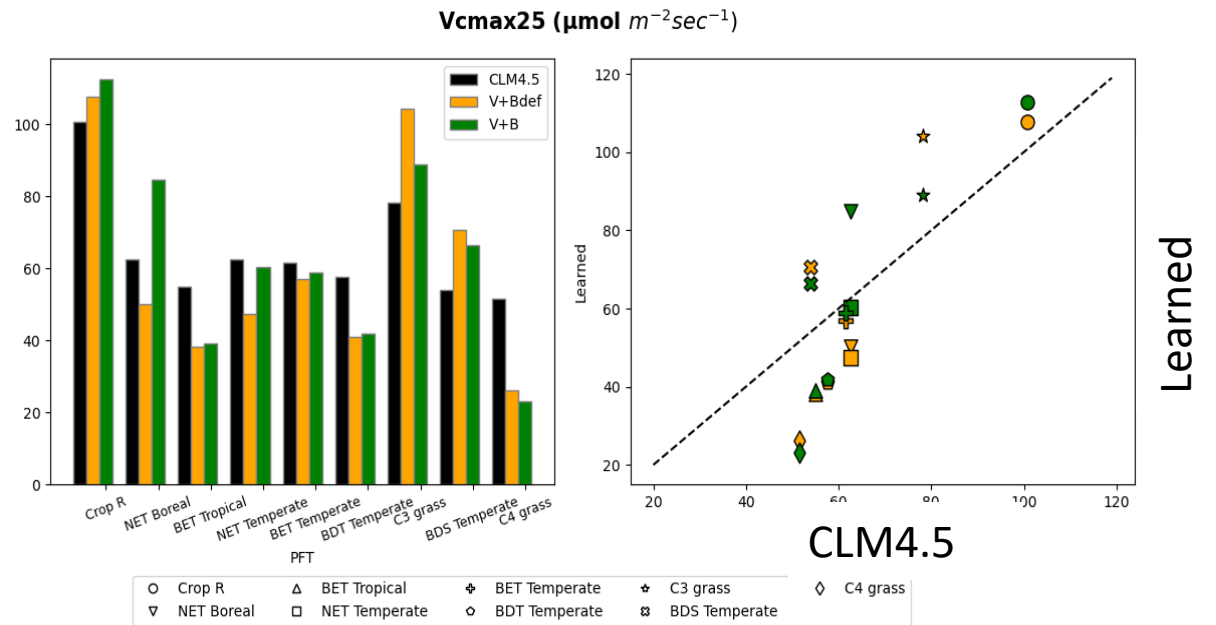
Runs	Corr		RMSE ( $\mu\text{mol m}^{-2} \text{s}^{-1}$ )		Bias ( $\mu\text{mol m}^{-2} \text{s}^{-1}$ )		NSE	
	Train	Test	Train	Test	Train	Test	Train	Test
<b>V<sub>def</sub>+B<sub>def</sub></b>	0.565		6.780		1.476		0.041	
V <sub>def</sub> +B <sub>def</sub> **	0.592		5.488		1.034		0.318	
V <sub>def</sub> +B	0.678	0.547	5.887	6.730	1.353	1.754	0.321	-0.084
V+B <sub>def</sub>	0.769	0.593	4.595	5.677	-0.129	-1.368	0.587	0.229
<b>V+B</b>	0.800	0.748	4.299	4.421	0.037	0.347	0.638	0.532
V+B**	0.774	0.768	4.269	4.198	0.056	0.092	0.597	0.581

\*\* refers to using C3\_only plants in dataset

Biogeosciences, 20, 2671–2692, 2023  
<https://doi.org/10.5194/bg-20-2671-2023>  
 © Author(s) 2023. This work is distributed under the Creative Commons Attribution 4.0 License.

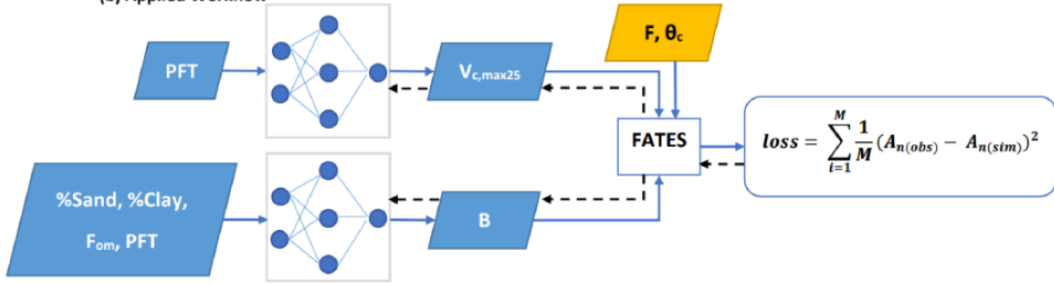
## A differentiable, physics-informed ecosystem modeling and learning framework for large-scale inverse problems: demonstration with photosynthesis simulations

Doaa Aboelyazeed<sup>1</sup>, Chonggang Xu<sup>2</sup>, Forrest M. Hoffman<sup>3,4</sup>, Jiangtao Liu<sup>1</sup>, Alex W. Jones<sup>5</sup>, Chris Rackauckas<sup>6</sup>, Kathryn Lawson<sup>1</sup> and Chuanming Shen<sup>1</sup>

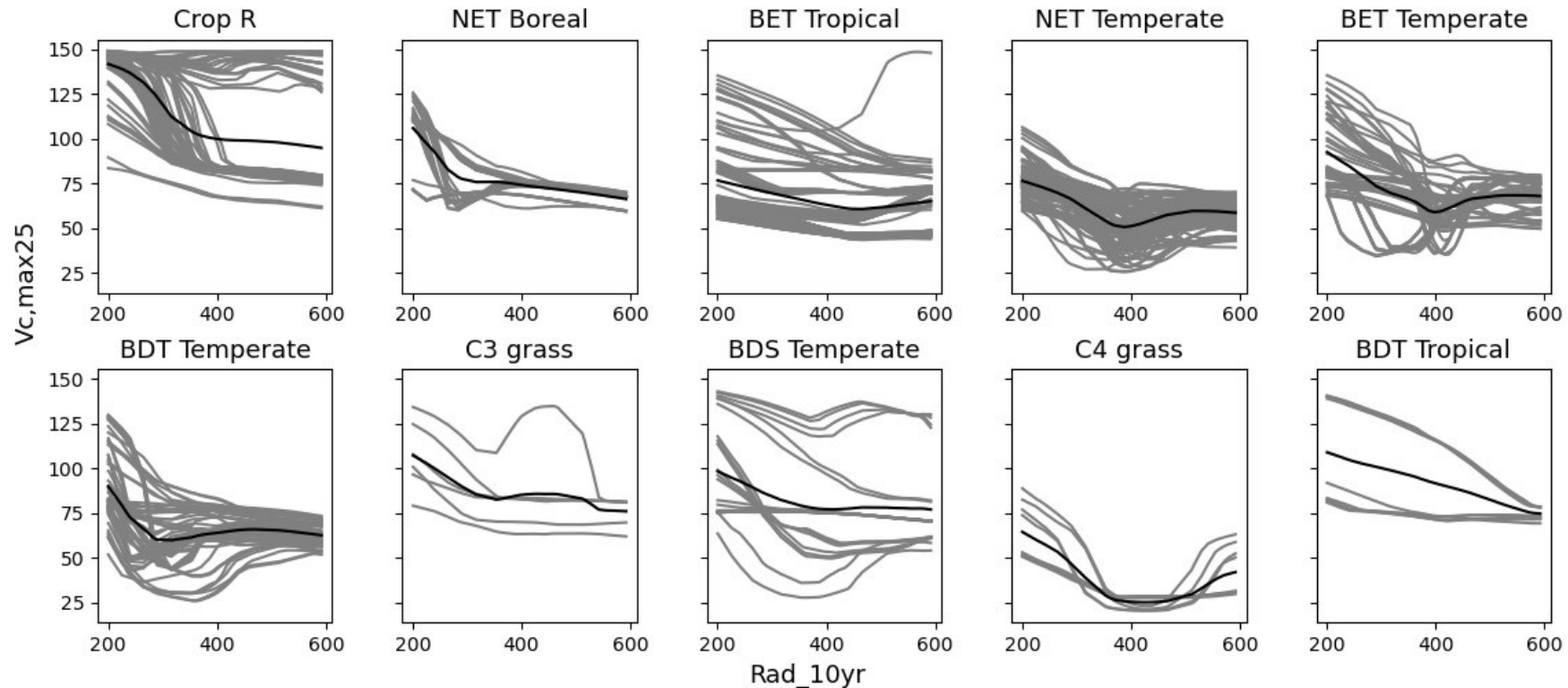


# Example 4. Ecosystem modeling: photosynthesis

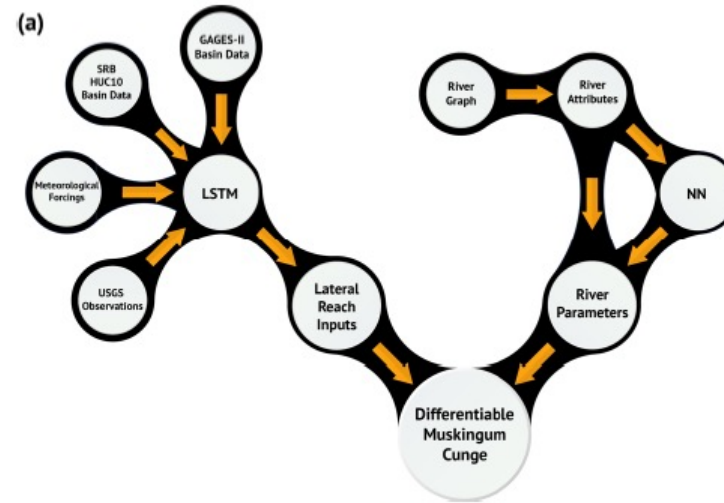
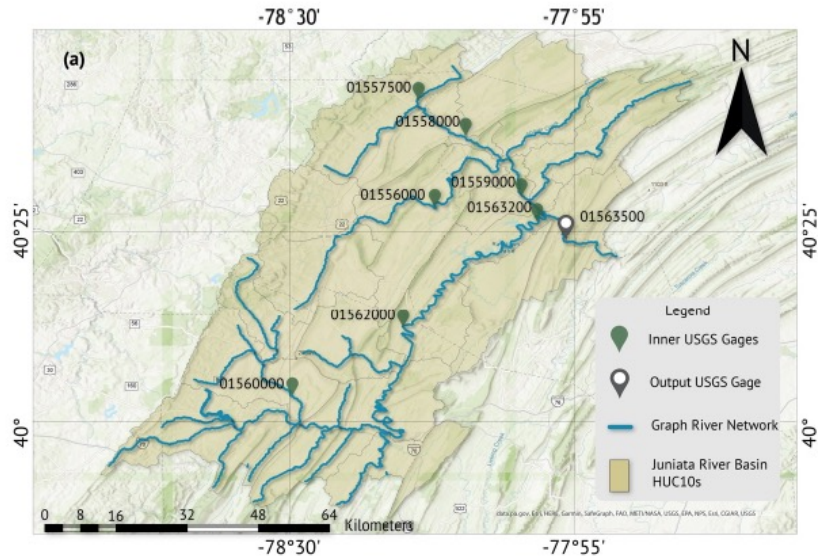
(b) Applied Workflow



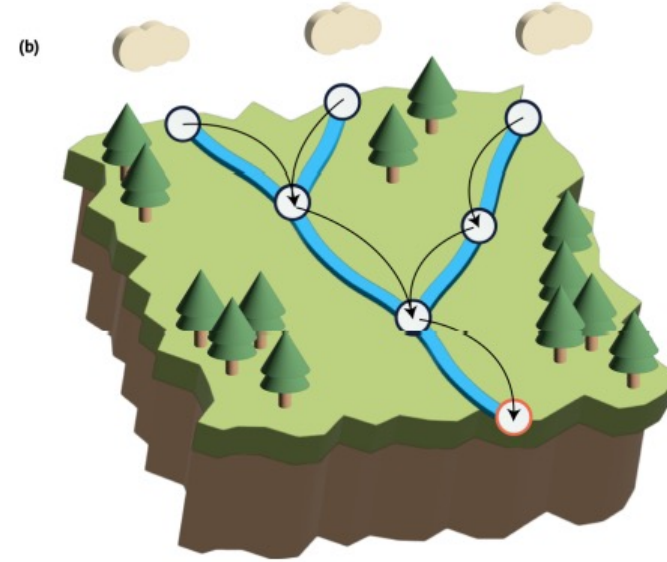
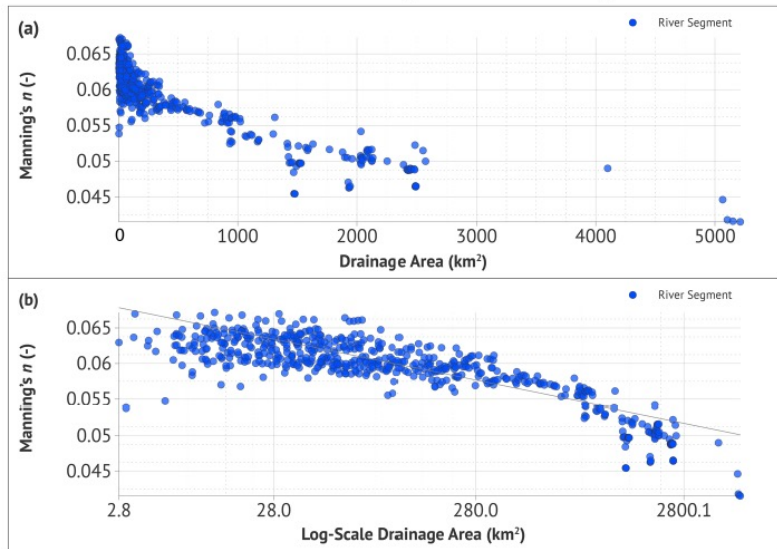
Discovering environmental dependencies of previously PFT-dependent parameter



# Example 4. Differentiable routing model



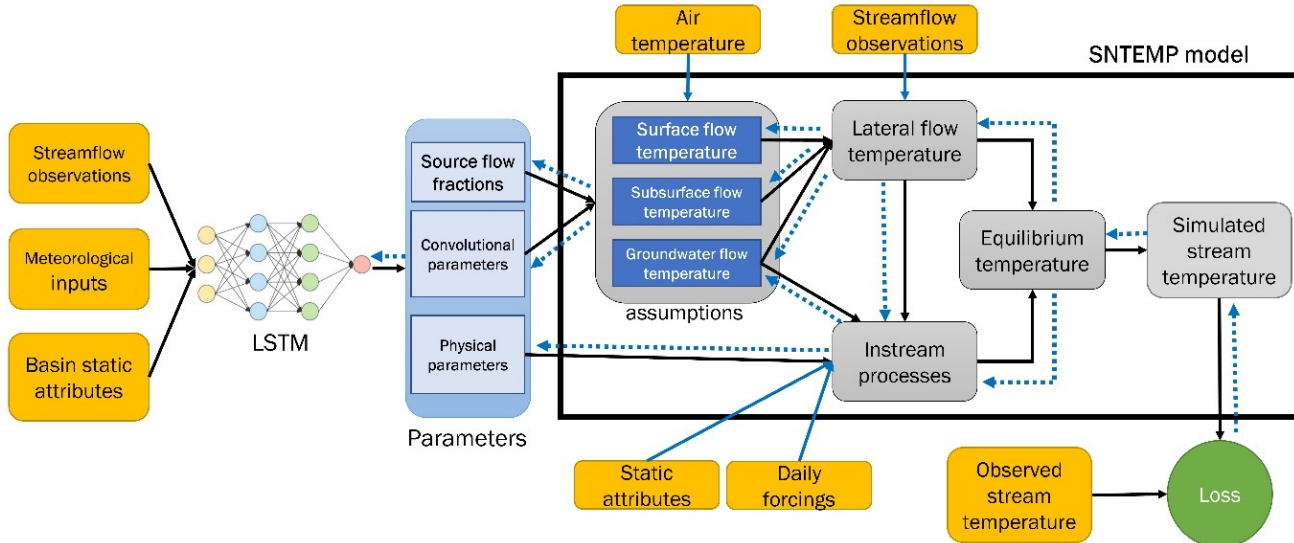
MLP  $n$  Distribution Trained Against Observed Discharge



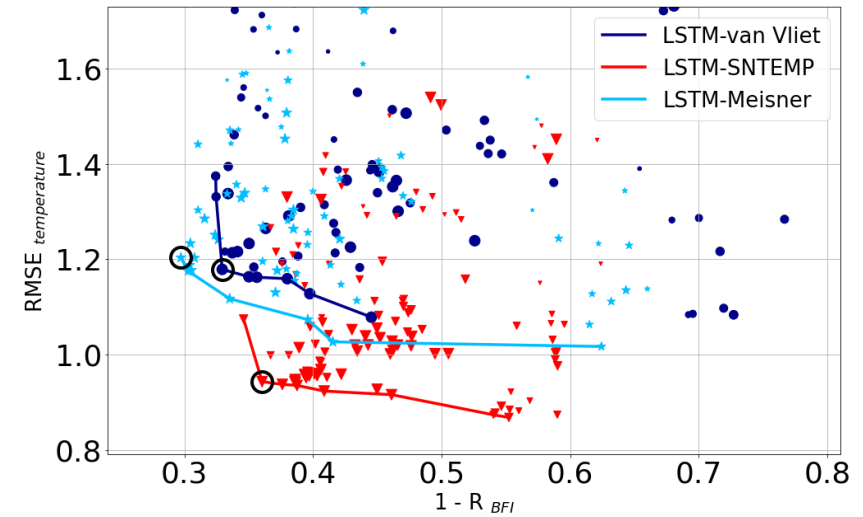


# Example 5. Water temperature modeling

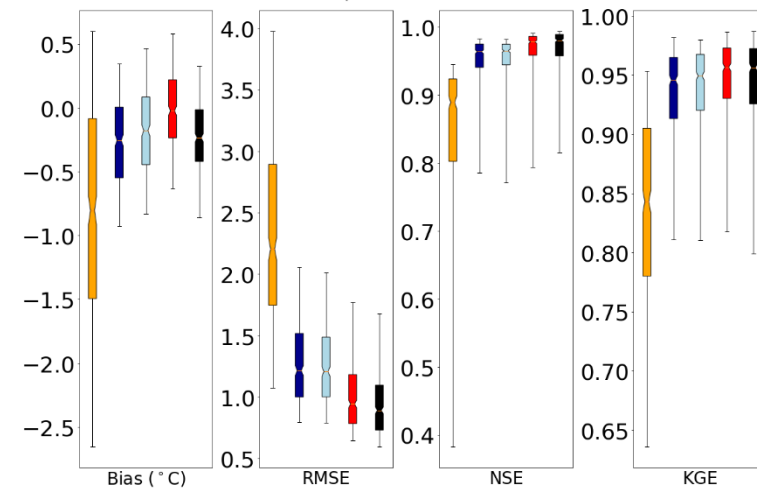
(a) The differentiable LSTM+SNTemp workflow



## Prior assumptions matter!

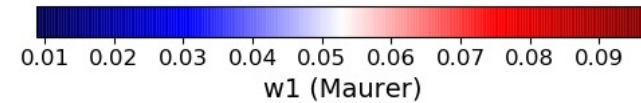
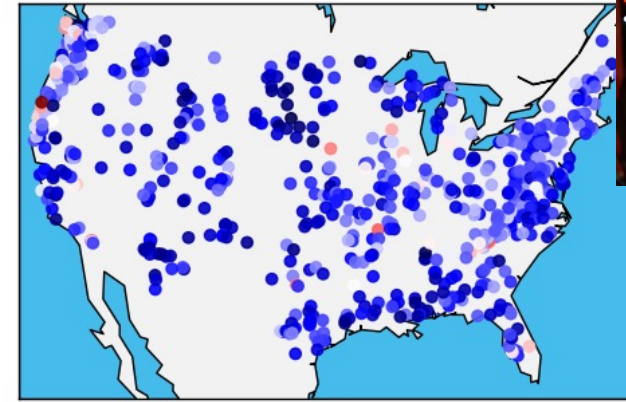
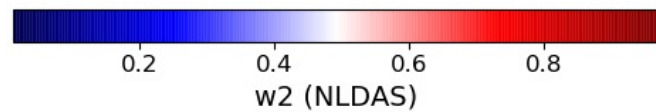
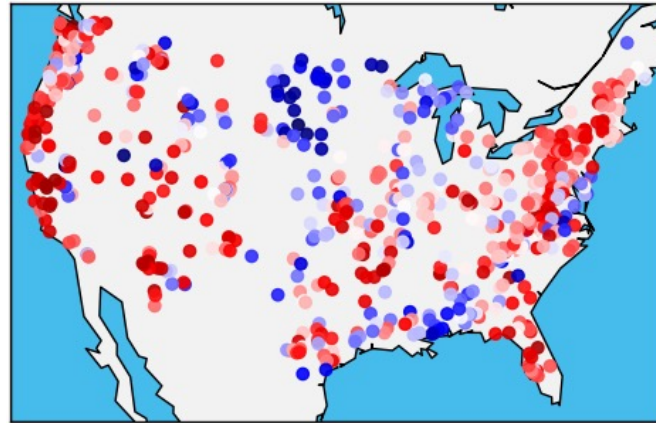
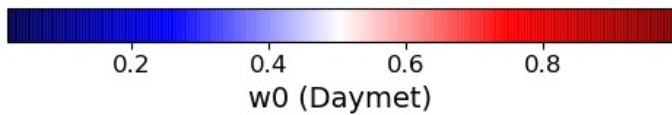
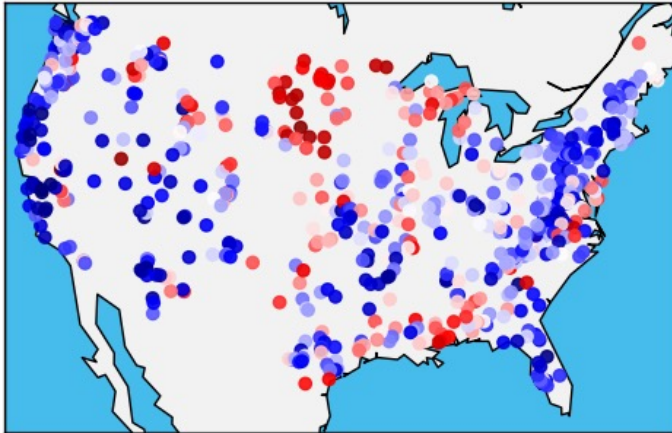


Comparing stream temperature performance of four models with previous studies



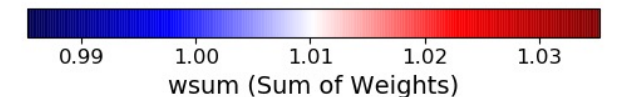
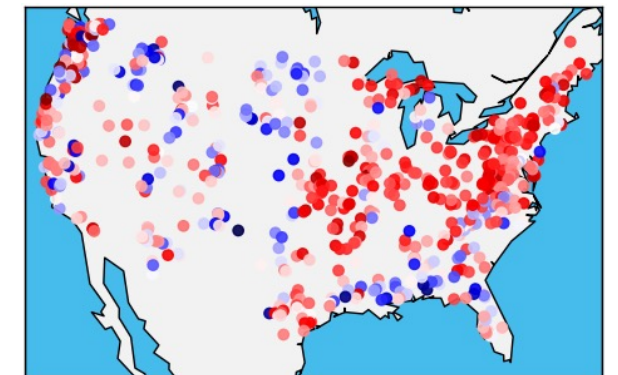
# Example 6. Fusion of forcings (in preparation)

NLDAS (0.56) > Daymet (0.41) > Maurer (0.03)



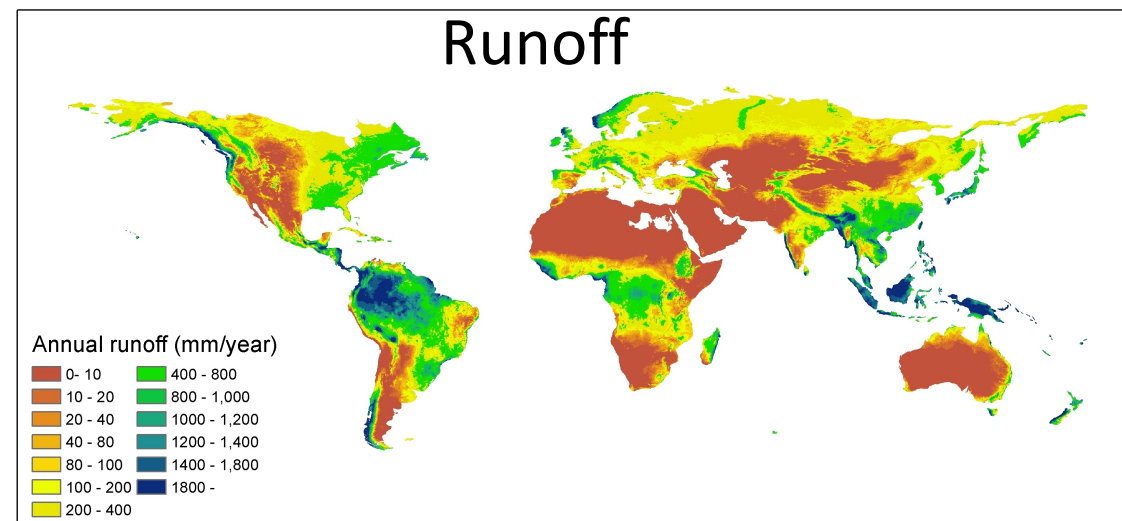
Simulation	Forcings	Median NSE	Median KGE	Low flow RMSE (mm/day)	High Flow RMSE (mm/day)
Single forcing w/o bias corr	Daymet	0.737	0.728	0.134	3.990
Multiforcing with bias correction	Daymet, Maurer, NLDAS	<b>0.770</b>	<b>0.780</b>	<b>0.082</b>	<b>3.414</b>

Low bias



# Future

- All kinds of models will be differentiable
- Climate change impact assessment will be done using high-quality models that have absorbed big data
- Many theories will be rewritten
- WaterGPT?





# Thank you!



@ChaopengShen  
cshen@engr.psu.edu

Hydroml.org

<https://github.com/mhpi>



Shen Multi-scale Hydrology, Processes and Intelligence Group (MHPI)

<http://water.engr.psu.edu/shen/hydroDL.html>

[CUAHSI cyberseminar series on BDML](#)

[WRR special issue on BDML](#)

[AGU Editor's review](#)

Hydrol. Earth Syst. Sci., 22, 5639–5656, 2018  
<https://doi.org/10.5194/hess-22-5639-2018>  
© Author(s) 2018. This work is distributed under the Creative Commons Attribution 4.0 License.



Hydrology and Earth System Sciences  
Open Access  
EGU

## HESS Opinions: Incubating deep-learning-powered hydrologic science advances as a community

Chaopeng Shen<sup>1</sup>, Eric Laloy<sup>2</sup>, Amin Elshorbagy<sup>3</sup>, Adrian Albert<sup>4</sup>, Jerad Bales<sup>5</sup>, Fi-John Chang<sup>6</sup>, Sangram Ganguly<sup>7</sup>, Kuo-Lin Hsu<sup>8</sup>, Daniel Kifer<sup>9</sup>, Zheng Fang<sup>10</sup>, Kuai Fang<sup>1</sup>, Dongfeng Li<sup>10</sup>, Xiaodong Li<sup>11</sup>, and Wen-Ping Tsai<sup>1</sup>

### Water Resources Research

#### REVIEW ARTICLE

10.1029/2018WR022643

#### Special Section:

Big Data & Machine Learning in Water Sciences: Recent Progress and Their Use in Advancing Science

#### A Transdisciplinary Review of Deep Learning Research and Its Relevance for Water Resources Scientists

Chaopeng Shen<sup>1</sup>

<sup>1</sup>Civil and Environmental Engineering, Pennsylvania State University, University Park, PA, USA

deepLDB

deepLDB -- a mac Landslide database

nature COMMUNICATIONS

ARTICLE

<https://doi.org/10.1038/s41467-021-26107-z> OPEN

From calibration to parameter learning: Harnessing the scaling effects of big data in geoscientific modeling

Wen-Ping Tsai<sup>1</sup>, Dapeng Feng<sup>1</sup>, Ming Pan<sup>2,3</sup>, Hylke Beck<sup>4</sup>, Kathryn Lawson<sup>1,5</sup>, Yuan Yang<sup>6,7</sup>, Jiangtao Liu<sup>1</sup> & Chaopeng Shen<sup>1,5</sup>✉