

A topographic map of the Eastern United States, showing the Appalachian Mountains and the Atlantic coast. The map is rendered in shades of green and brown, with blue lines indicating rivers and water bodies. The text is overlaid on the left side of the map.

# Network Analysis: Data-Driven Optimization

## Tools & Insights

Matthew Cashman

Leah Staub, Joel Blomquist,

Zach Clifton, John Hammond

MD-DE-DC Water Science Center  
US Geological Survey

This information is preliminary and is subject to revision. It is being provided to meet the need for timely best science. The information is provided on the condition that neither the U.S. Geological Survey nor the U.S. Government shall be held liable for any damages resulting from the authorized or unauthorized use of the information.

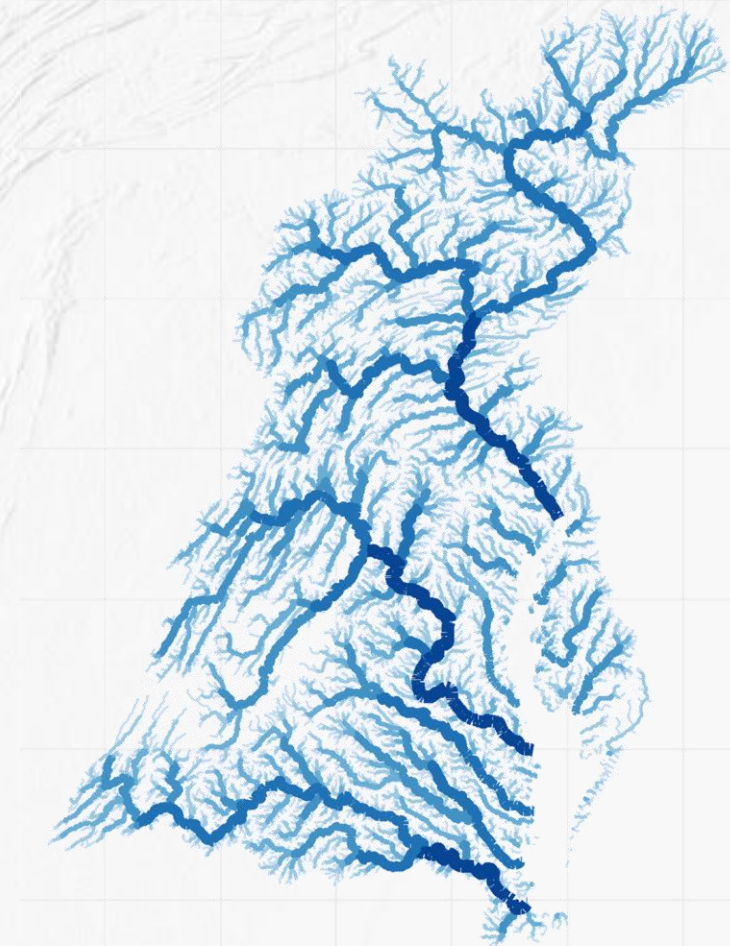
# Outline

- Introduction
- A “Reusable” Workflow
- Example USGS Use-Cases
- Quantifying Gage Loss



# Network Analysis

- A comparison of the representativeness of a monitoring network (“sample” or “target”) versus the underlying population it is supposed to represent (e.g. the river network)
- Useful for:
  - Predictive modeling
  - Identifying monitoring gaps to fill
  - Minimizing the impacts of monitoring cuts
  - Others...?

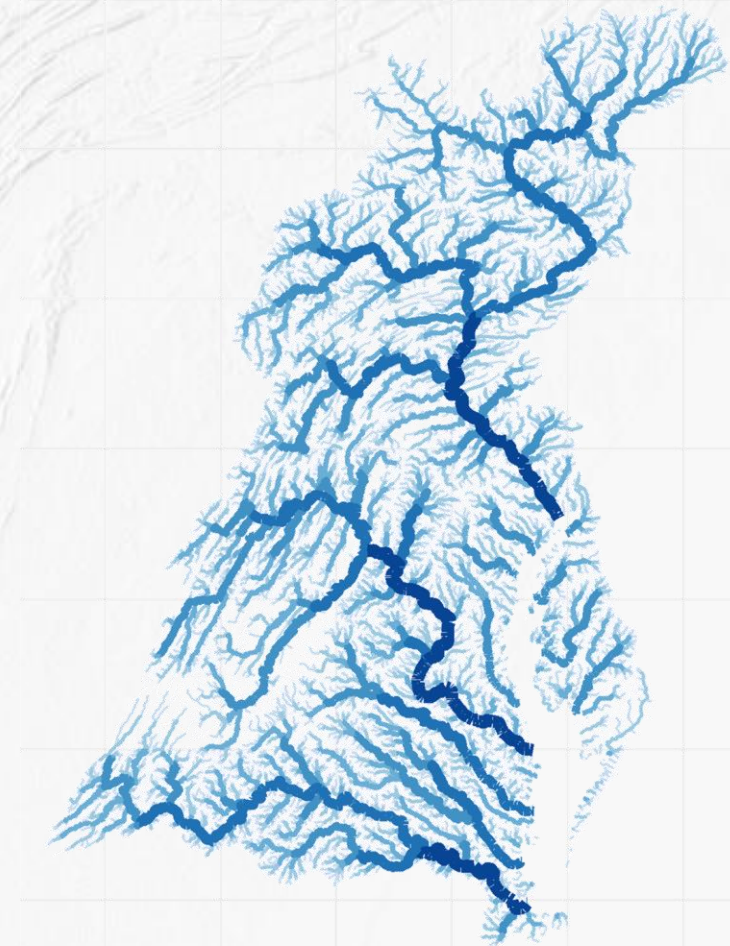


# Network Analysis

A basic question: “Are the samples representative of the population?”

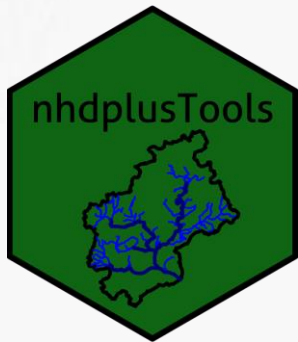
- But challenges:
  - Unwieldy “big” data (large sizes, different sources)
  - Large GIS preprocessing \*before\* any analysis
  - Analysis methods (“how”)
  - Method reproducibility and iterative changes

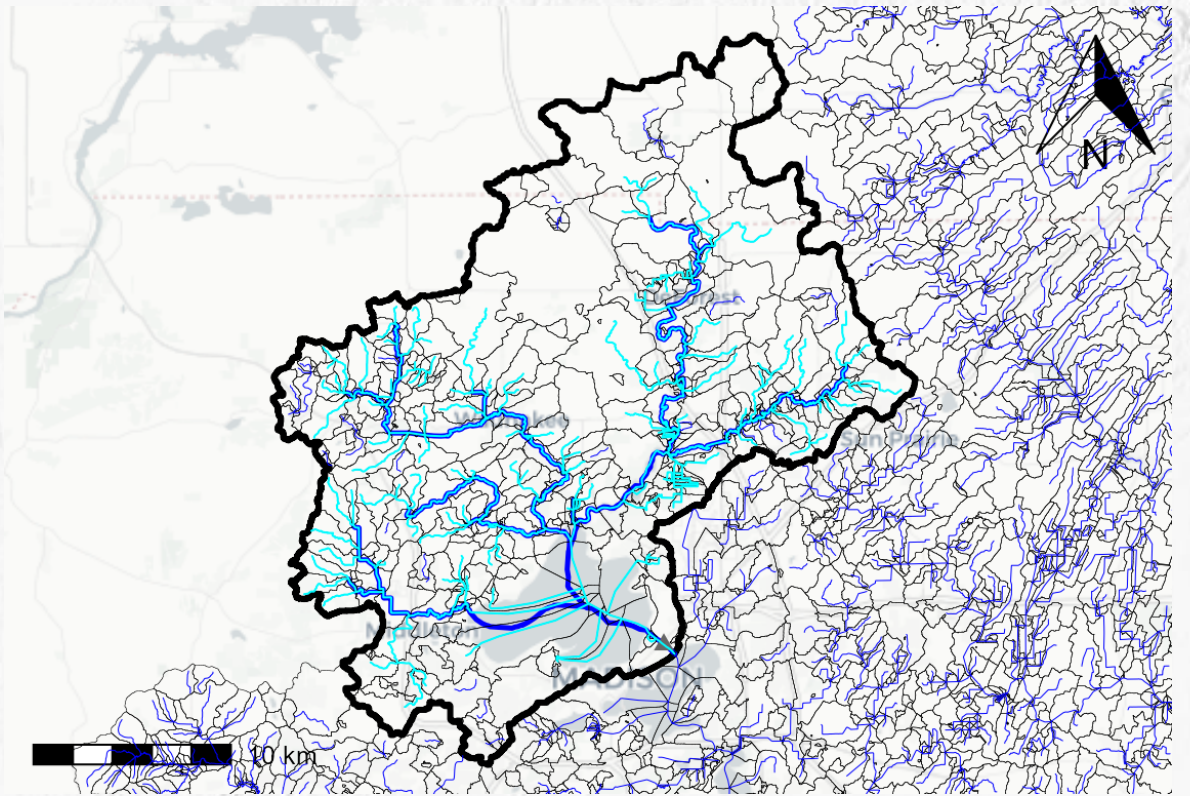
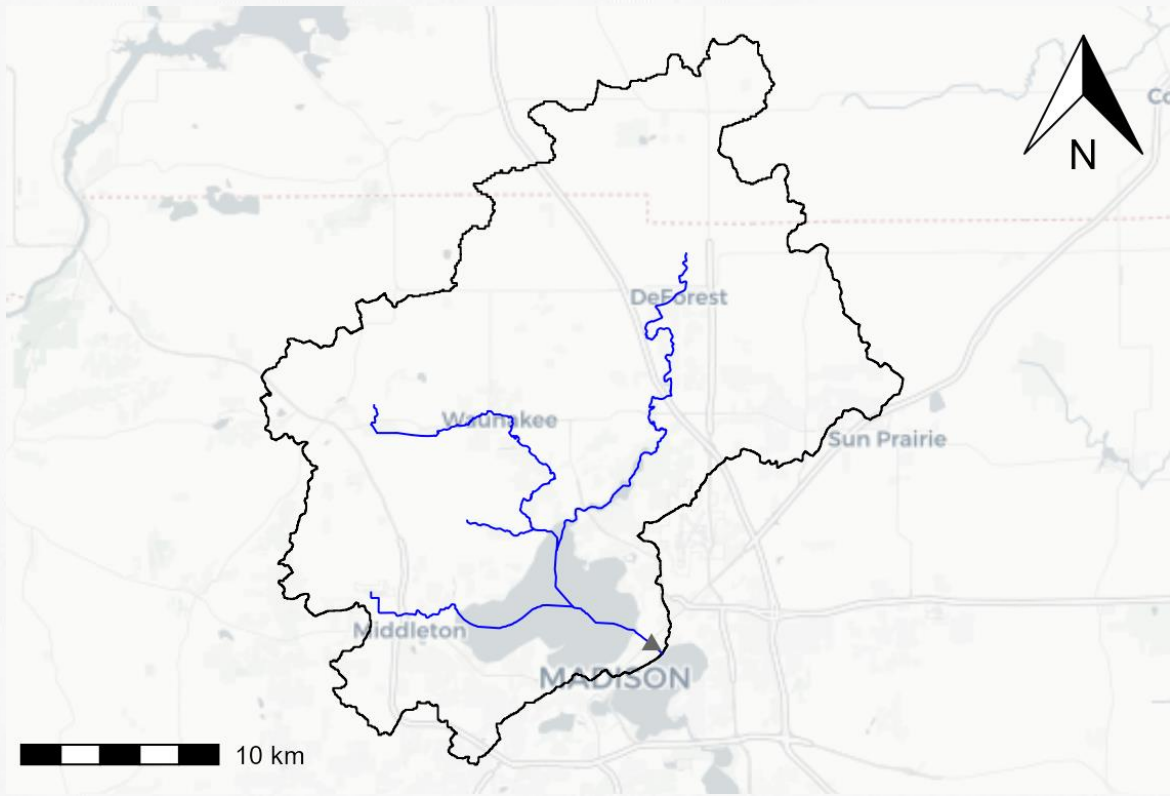
➤ “Representative of *what?*”



# Making Network Analysis “easier”

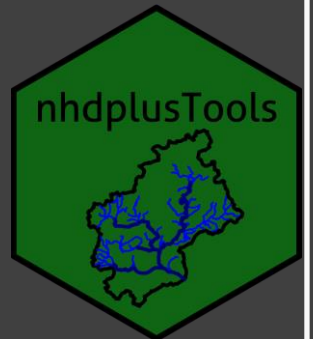
- *In-development* R code repository at [www.code.usgs.gov](http://www.code.usgs.gov)
- Tools to streamline network analyses (some R skills involved)
  - Setup and download data
  - Custom queries (landscape variables and monitoring locations)
  - Reproducible for different networks/focal areas
- Leverages David Blodgett’s nhdplusTools package





## Using NHDPlus v2.1 1:100k Framework

- Using linked data framework
  - Flowlines
  - Catchments



# Gathering required data

1. NHDPlus v2.1 CONUS seamless dataset directly from EPA<sup>1</sup>
  - *get\_nhdplus2.1\_seamless\_l48()*
2. 1,000s of upstream variables for 2.7 mil catchments (CONUS) from USGS<sup>2</sup>
  - *get\_mikew\_nhdplus\_accum\_attributes()*
3. Reach COMID/HUC crosswalk table from USGS<sup>3</sup>
  - *get\_comid\_crosswalk()*

• <sup>1</sup><https://www.epa.gov/waterdata/nhdplus-national-data>

• <sup>2</sup>Wieczorek, M.E., Jackson, S.E., and Schwarz, G.E., 2018, Select Attributes for NHDPlus Version 2.1 Reach Catchments and Modified Network Routed Upstream Watersheds for the Conterminous United States (ver. 3.0, January 2021): U.S. Geological Survey data release, <https://doi.org/10.5066/F7765D7V>.  
<https://www.sciencebase.gov/catalog/item/5669a79ee4b08895842a1d47>

• <sup>3</sup>Moore, R.B., Johnston, C.M., and Hayes, L., 2019, Crosswalk Table Between NHDPlus V2.1 and its Accompanying WBD Snapshot of 12-Digit Hydrologic Units: U.S. Geological Survey data release, <https://doi.org/10.5066/P9CFXHGT>. <https://www.sciencebase.gov/catalog/item/5c86a747e4b09388244b3da1>

# Create local GIS river network

- Subset NHDPlus v2.1 to area of interest and save as smaller geodatabase for future use
  - subsetNHD()
    - Takes any HUC combination: HUC02 → HUC12
    - Single or in combination
  - Run once
    - Repeat analysis
  - Avoids RAM limitations via SQL queries



# Landscape queries

17.0 GB  
134 folders  
772 files

1. **Best Management Practices** - characteristics such as agricultural management practices and land in conservation practices.
2. **Chemical** - characteristics such as nitrogen application or toxicity weighted use.
3. **Climate and Water Balance Model** - characteristics such as model outputs of runoff, actual evapotranspiration or ground water storage.
4. **Climate** - characteristics such as mean precipitation, temperature, relative humidity, or evapotranspiration.
5. **Geology** - characteristics such as Hunt or Soller surficial geologies.
6. **Hydrologic** - characteristics such as base flow or infiltration excess overland flow. Hydrologic Modifications, characteristics such as dam storage or tile drains.
7. **Hydrologic Modifications** - characteristics such as dam storage or tile drains.
8. **Landscape** - characteristics such as NLCD, CDL or NWALT.
9. **Population Infrastructure** - characteristics such as population, housing, and road densities.
10. **Regions** - characteristics such as EcoRegions, Physiography or Hydrologic Landscapes.
11. **Soils** - characteristics such as STATSGO, soil salinity, and soil restrictive layer.
12. **Topographic Characteristics** - characteristics such as basin area, slope and elevation.
13. **Water use** - characteristics such as estimated freshwater withdrawals and estimated freshwater consumption by thermo-electric power plants

.zip

# Landscape queries

## what\_attributes()

- Best\_Management
- Chemical\_Attributes
- Climate\_Attributes
- Climate\_Water
- Geologic\_Attributes
- Hydrologic\_Attributes
- Hydrologic\_Modification
- Land\_Cover
- Population\_Infrastructure
- Regional\_Attributes
- Soil\_Attributes
- Topographic\_Attributes
- Water\_Use

# Landscape queries

what\_attributes(directory = "Climate\_Water")

## Climate\_Water

- Annual\_Average\_Actual\_Evapotranspiration\_millimeters\_from\_2014\_\_2015
- Annual\_Average\_Potential\_Evapotranspiration\_millimeters\_from\_2014\_\_2015
- Annual\_Average\_Precipitation\_millimeters\_from\_1945\_2015
- Annual\_Average\_Runoff\_millimeters\_from\_1945\_\_2015
- Annual\_Average\_Temperature\_Celsius\_from\_1945\_\_2015
- an\_Annual\_Water\_Balance\_Variables\_for\_the\_Period\_of\_Record\_2000\_2014\_and\_Detrended\_for\_the\_Year\_2012
- Average\_Annual\_Water\_Balance\_Variables\_over\_the\_Period\_of\_Record\_2000\_2014
- Average\_Monthly\_Runoff\_mm\_from\_McCabe\_and\_Wolock\_s\_Runoff\_Model\_Over\_Period\_of\_Record\_1951\_2000
- Monthly\_Average\_Annual\_Runoff\_mm\_from\_McCabe\_and\_Wolock\_s\_Runoff\_Model\_1945\_\_2015
- Monthly\_Average\_Precipitation\_millimeters\_from\_1945\_\_2015
- Monthly\_Average\_Temperature\_Celsius\_from\_1945\_\_2015

# Landscape queries

what\_attributes(directory = "Geologic\_Attributes")

## Geologic\_Attributes

- Attributes\_for\_NHDPlus\_Version\_2\_1\_Flowlines\_for\_the\_Conterminous\_United\_States\_Lithology
- Bedrock\_Permeability\_Classes
- Combination\_of\_landuse\_and\_geology\_
- Generalized\_Geology\_Type\_Reed\_and\_Bush\_2001\_
- Hunt\_Geology\_1999
- Principal\_Aquifers\_and\_Rock\_Types
- Select\_Geochemical\_Characteristics\_Based\_on\_Olson\_Geology\_Types
- Soller\_Surficial\_Materials\_2009

# Landscape queries

what\_attributes(directory = “Hydrologic\_Modification\_Attributes”)

## Hydrologic\_Modification

- Attributes\_for\_Surface\_Water\_Impoundments
- Estimated\_Area\_of\_Agricultural\_Land\_Drained\_by\_Field\_Ditches\_1992
- Estimated\_Area\_of\_Subsurface\_Drainage\_on\_Agricultural\_Land\_1992
- Estimates\_of\_Subsurface\_Tile\_Drainage\_Extent\_for\_the\_Conterminous\_United\_States\_Early\_1990s
- Major\_Sites\_of\_the\_National\_Pollutant\_Discharge\_Elimination\_System\_NPDES\_
- National\_Inventory\_of\_Dams\_NID\_Storage\_and\_Construction\_by\_Decade\_1930\_to\_2010
- Percent\_of\_Irrigated\_Agriculture\_2002\_MlrAD\_data
- Percent\_of\_Irrigated\_Agriculture\_2007\_MlrAD\_data
- Percent\_of\_Irrigated\_Agriculture\_2012\_MlrAD\_data

# Landscape queries

what\_attributes(directory = "Land\_Cover")

## Land\_Cover

- 2012\_Conservation\_Reserve\_Program\_land\_allocated\_to\_2012\_NWALT\_landuse
- Crop\_Land\_Data\_Layer\_2012
- Crop\_Land\_Data\_Layer\_2013
- Crop\_Land\_Data\_Layer\_2014
- the\_conterminous\_United\_States\_seasonal\_enhanced\_vegetation\_index\_fall\_of\_2011\_through\_summer\_of\_2012
- Modeled\_historical\_land\_use\_and\_land\_cover\_for\_the\_conterminous\_United\_States\_1992\_2002
- National\_Land\_Cover\_Database\_2001\_NLCD\_2001\_
- National\_Land\_Cover\_Database\_2006\_NLCD\_2006\_
- National\_Land\_Cover\_Database\_2011\_NLCD\_2011\_
- National\_Land\_Cover\_Database\_2016\_Versions\_for\_the\_Years\_2001\_2004\_2006\_2008\_2011\_2013\_and\_2016
- NAWQA\_Wall\_to\_Wall\_Anthropogenic\_Land\_Use\_Trends\_NWALT\_Timber\_1999\_2012
- NLCD\_2001\_Percent\_Imperviousness
- NLCD\_2001\_Percent\_Imperviousness\_in\_100\_Meter\_Riparian\_Buffer
- NLCD\_2006\_Percent\_Imperviousness
- NLCD\_2006\_Percent\_Imperviousness\_in\_100\_Meter\_Riparian\_Buffer
- NLCD\_2011\_Percent\_Imperviousness
- Percent\_NLCD\_2011\_in\_50\_Meter\_Riparian\_Buffer
- Percent\_NLCD\_2011\_Tree\_Canopy
- Percent\_of\_Watershed\_Covered\_by\_Waterbodies\_from\_NHD\_High\_Resolution
- Wildfire\_2000\_2012

# Landscape queries

`what_attributes`(Subdirectory: National\_Land\_Cover\_Database\_2016\_re-release)

- NHDV2\_NLCD2016s\_Revised.xml
- NLCD01\_ACC\_CONUS.zip
- NLCD01\_CAT\_CONUS.zip
- NLCD01\_TOT\_CONUS.zip
- NLCD04\_ACC\_CONUS.zip
- NLCD04\_CAT\_CONUS.zip
- NLCD04\_TOT\_CONUS.zip
- NLCD06\_ACC\_CONUS.zip
- NLCD06\_CAT\_CONUS.zip
- NLCD06\_TOT\_CONUS.zip
- NLCD08\_ACC\_CONUS.zip
- NLCD08\_CAT\_CONUS.zip
- NLCD08\_TOT\_CONUS.zip
- NLCD11\_ACC\_CONUS.zip
- NLCD11\_CAT\_CONUS.zip
- NLCD11\_TOT\_CONUS.zip
- NLCD13\_ACC\_CONUS.zip
- NLCD13\_CAT\_CONUS.zip
- NLCD13\_TOT\_CONUS.zip
- NLCD16\_ACC\_CONUS.zip
- NLCD16\_CAT\_CONUS.zip
- NLCD16\_TOT\_CONUS.zip
- Version\_History.txt

# Batch load directly from zips

`load_attributes`(Subdirectory: National\_Land\_Cover\_Database\_2016\_re-release)

- Batch import of subdirectory of attributes
- Reads zip files directly
- Fast
  - Benchmark: 2 min for ~1billion NLCD values (2.8 million records x 133 variables)
- Delivers subset table for Area of Interest
- Saves CSV to disk for future easy reference

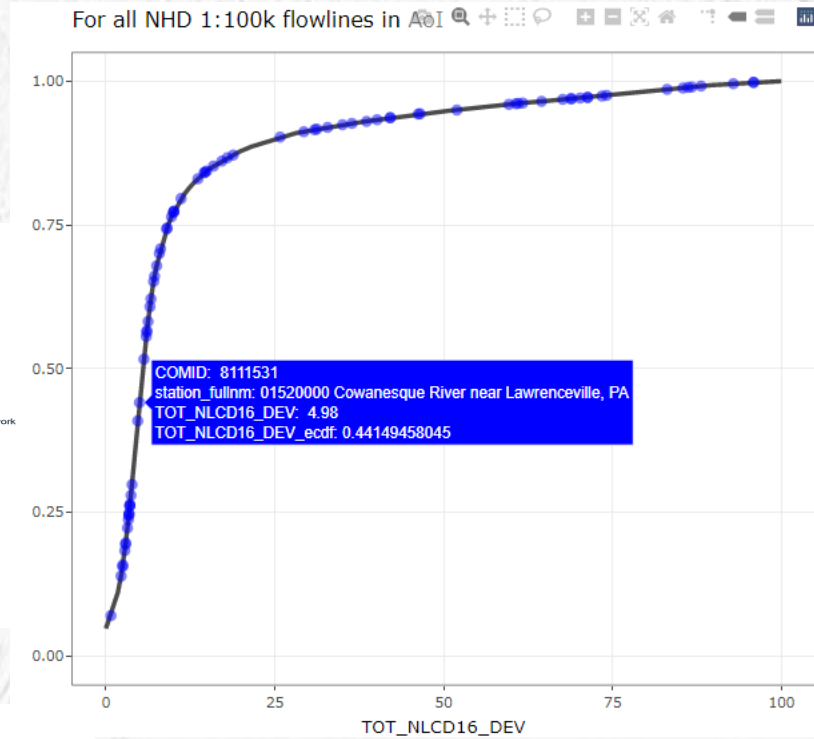
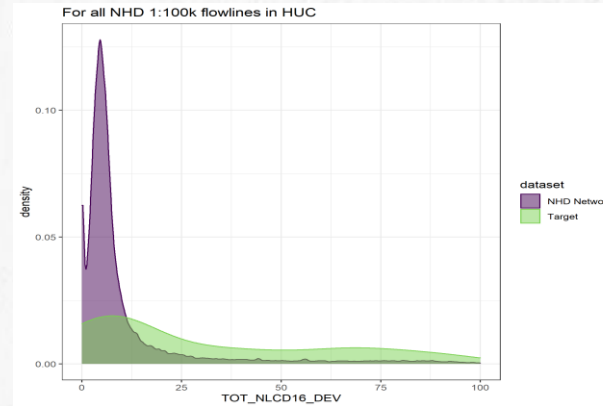
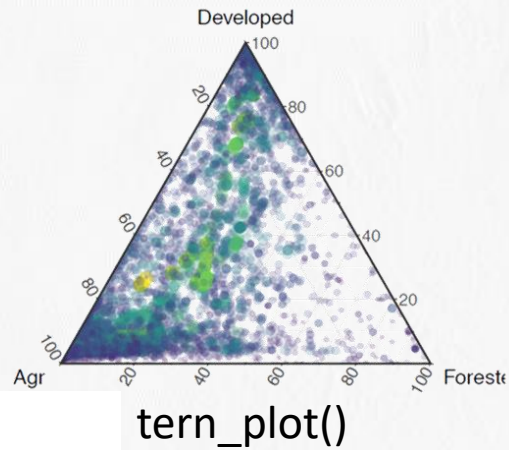


# Gage queries

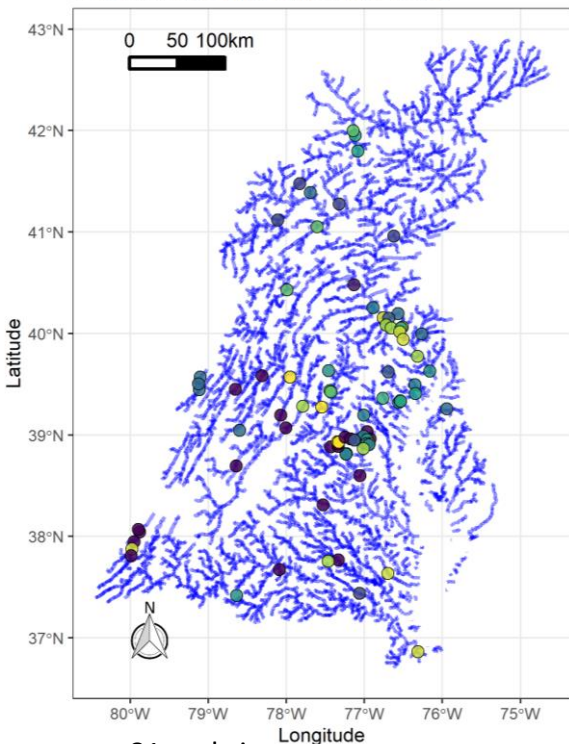
## GetNWISGages()

- HUCs: single or combination
- Parameters: USGS code (e.g. 00010: Water Temperature)
- Data type: Real-time or discrete
- Active: Currently active or historical

# Static & Interactive Plotting



NHD Flowlines for Area of Interest  
Location of Real-Time Temperature Gages

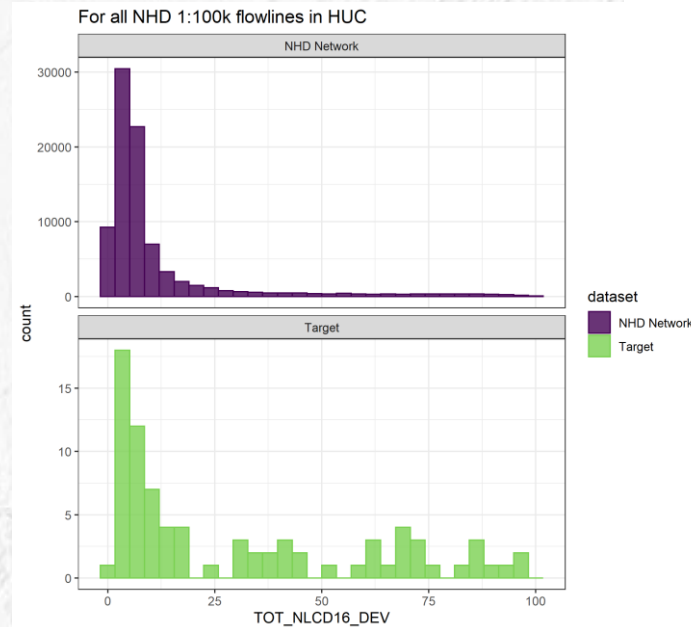


log10(TotDASqKM + 1) · 0 · 1 · 2 · 3 · 4

RT Data  
Start Year  
2020  
2016  
2012  
2008

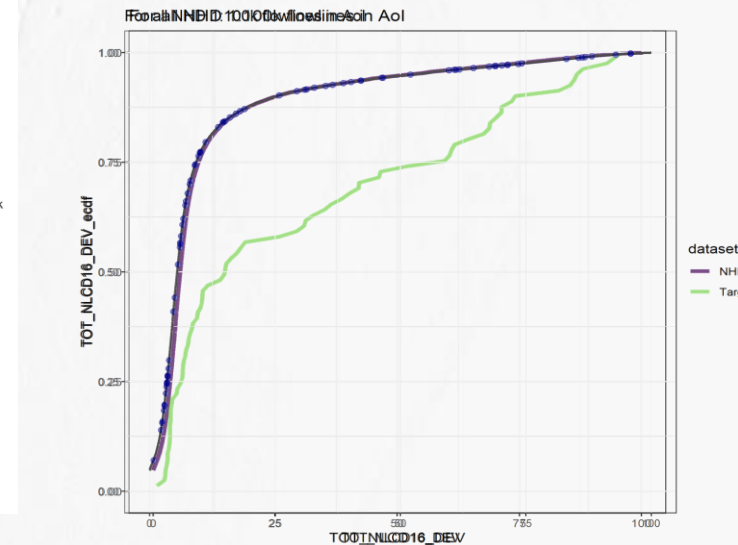
81 real-time water temperature gages

Density\_analysis()



Histogram\_analysis()

Distribution\_analysis()



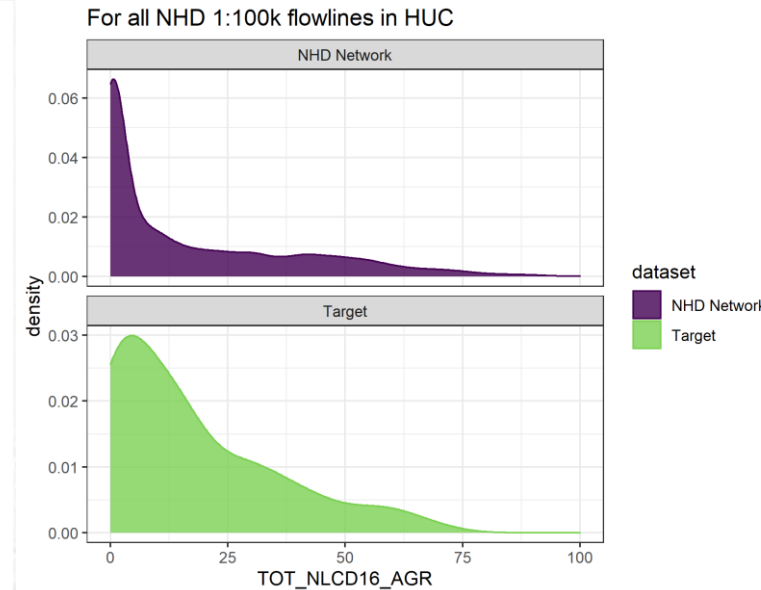
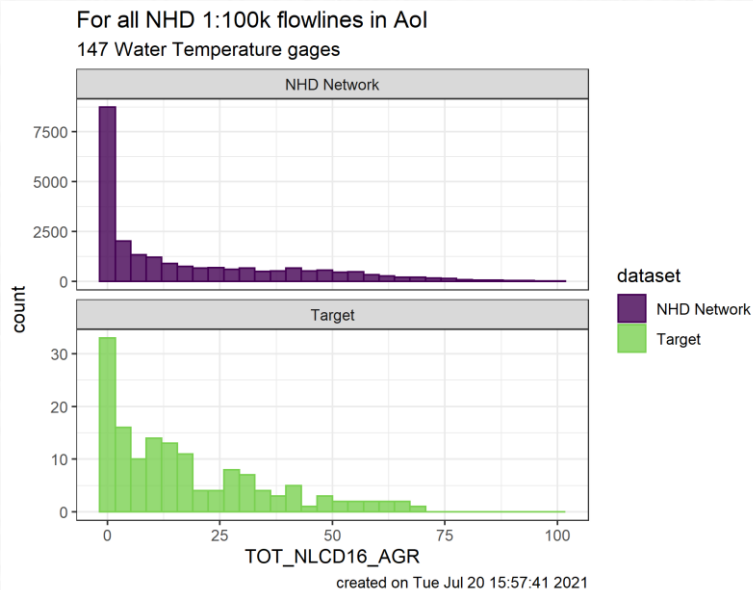
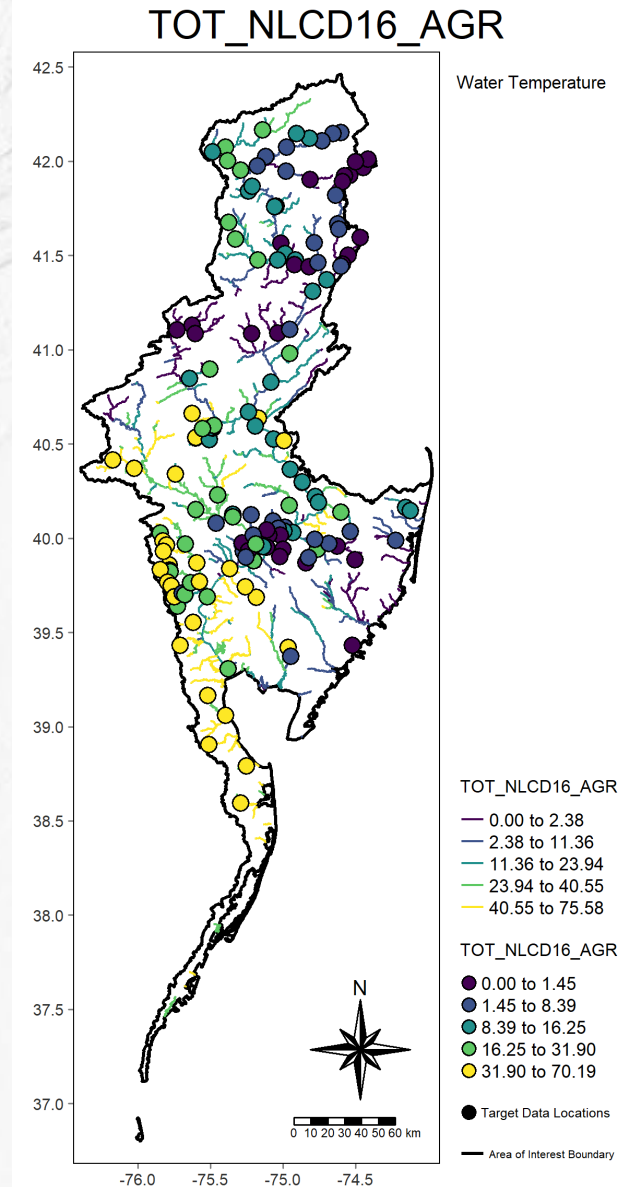
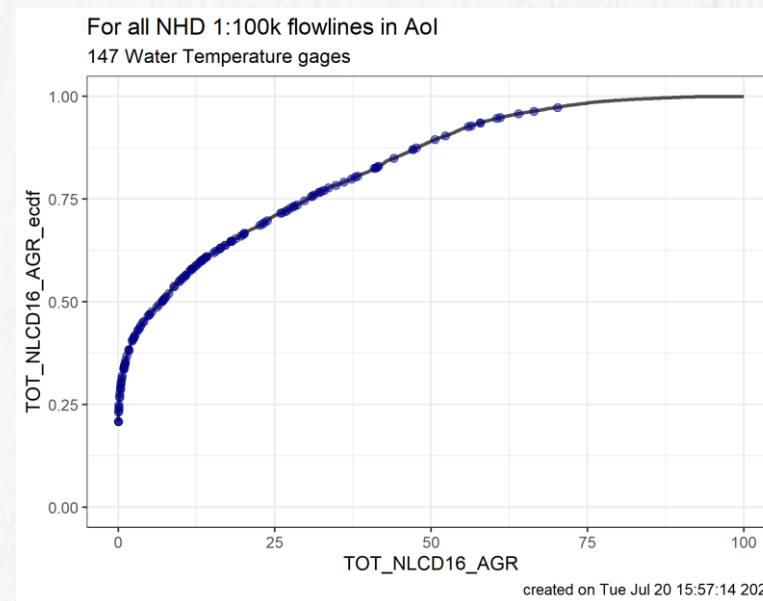
Data are provisional and subject to revision

# Continued analysis

- Data are all in R for further or custom analysis
  - Spatial data in *sf* (simple features) format
- Export tables to CSV
- Also easy export to ESRI/Arc friendly formats

# Use Case: Next Generation Water Observation System

- Delaware River Basin
- Objective: Temperature
  - Assess before/after 2019 gage additions

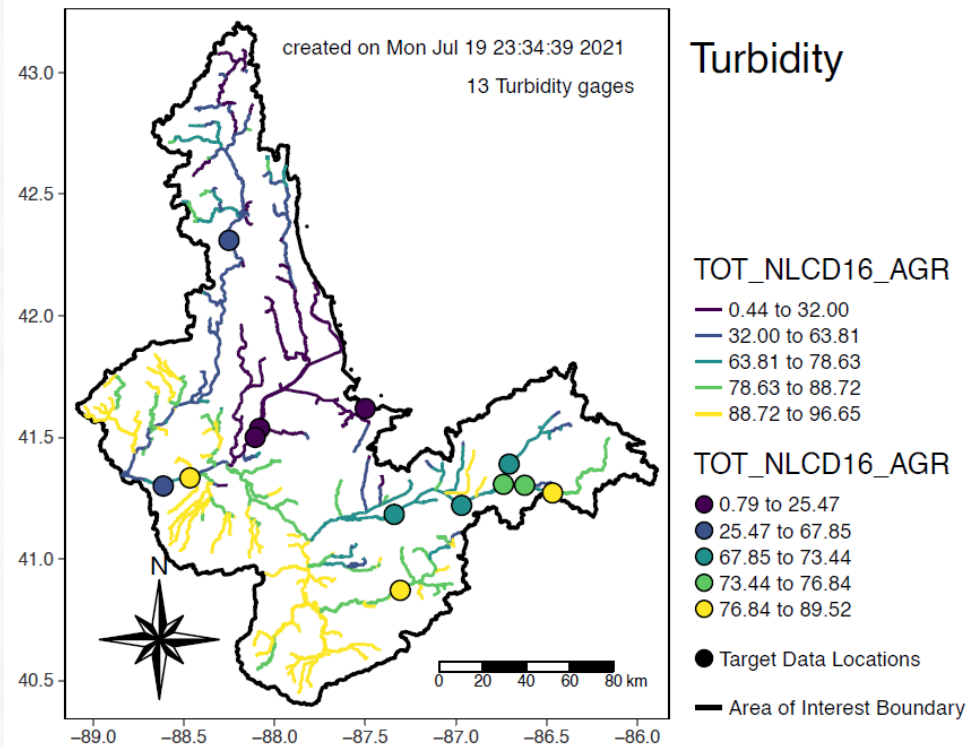


Data are provisional and subject to revision

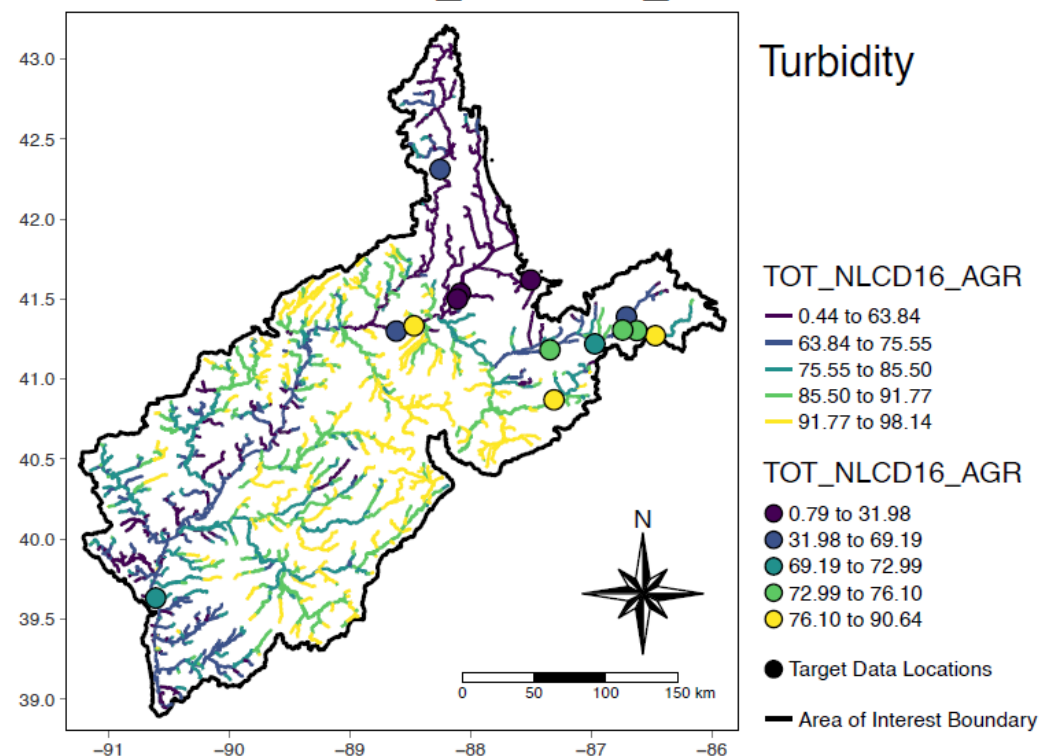
# Use Case: Next Generation Water Observation System

- Upper Illinois River
- Objective: Temp, SpC, DO, pH, Turbidity

## Upper Illinois River TOT\_NLCD16\_AGR



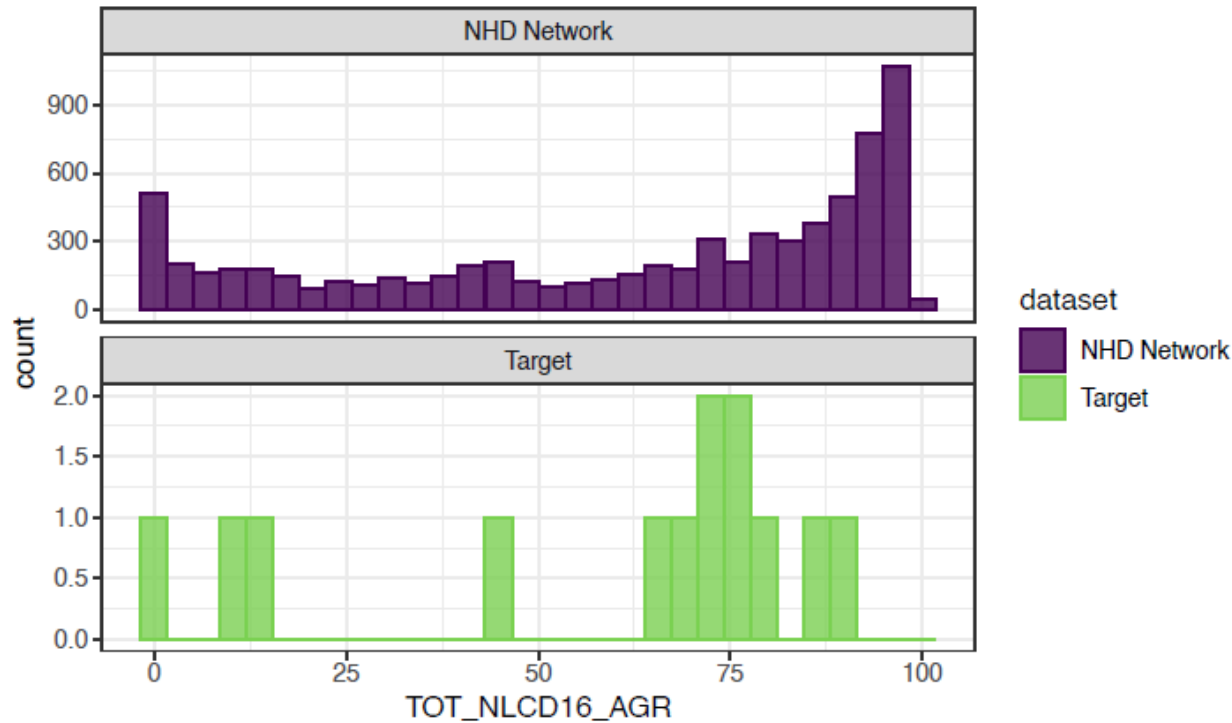
## Illinois River TOT\_NLCD16\_AGR



# Use Case: Next Generation Water Observation System

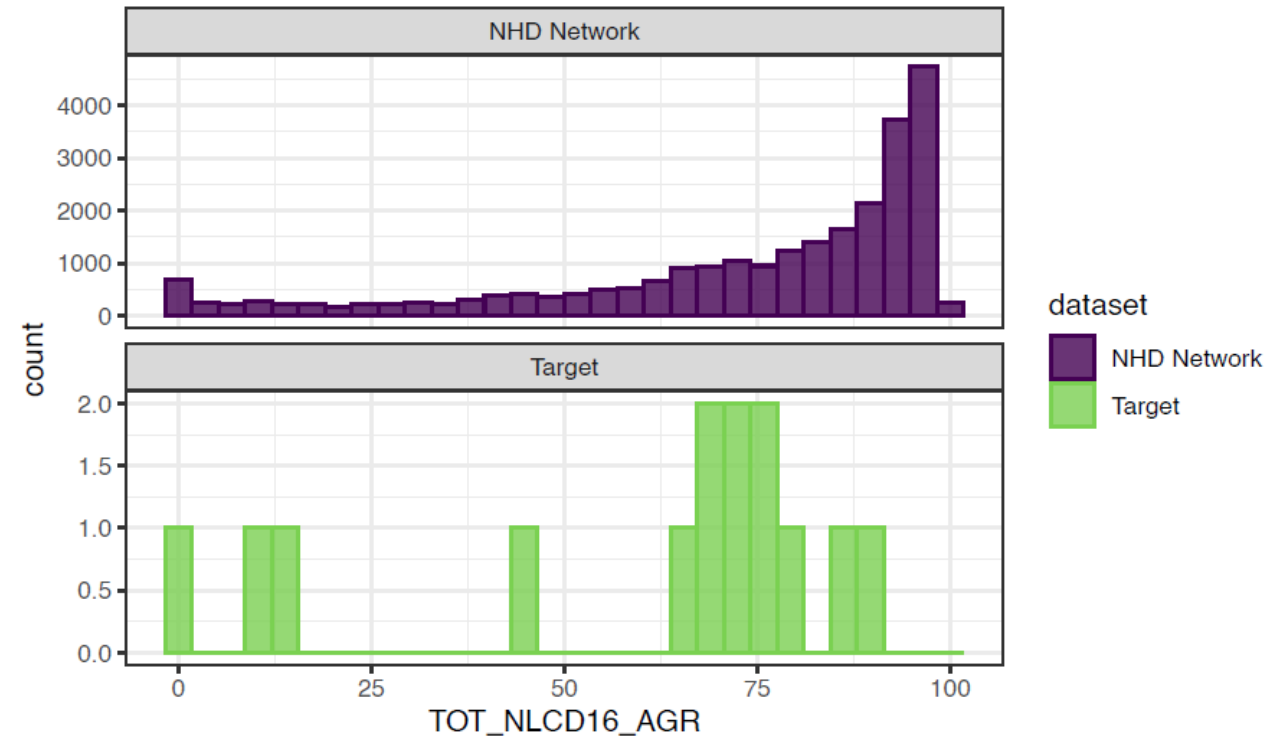
## Upper Illinois River

For all NHD 1:100k flowlines in Aol  
13 Turbidity gages



## Illinois River

For all NHD 1:100k flowlines in Aol  
14 Turbidity gages



# Making analyses actionable

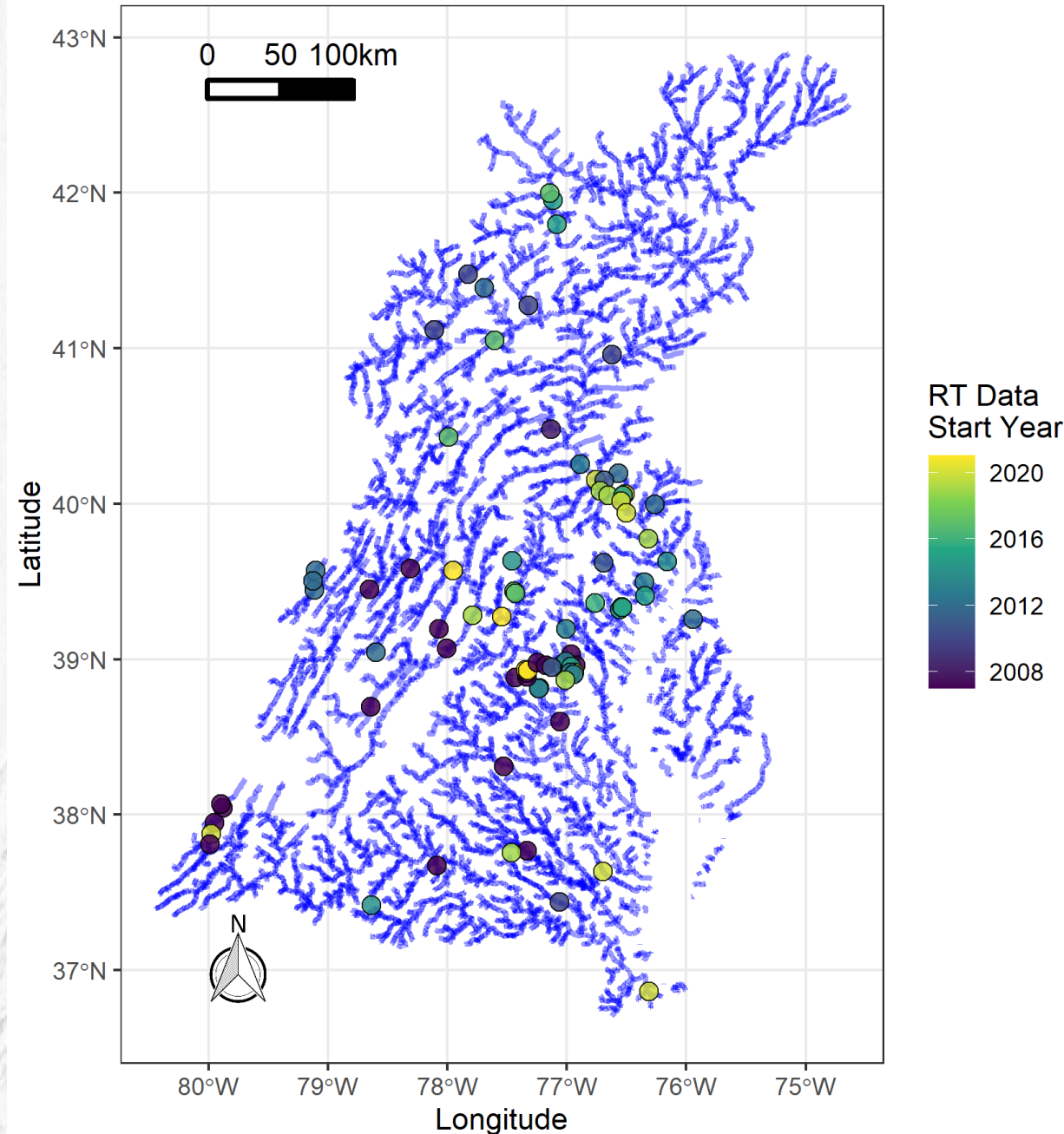
- These analyses so far have been largely qualitative:
  - Different graphical ways to compare networks and variable distributions
- How do you use these for decision-making?
  1. Gap identification and filling (network enhancements - adding)
  2. Quantifying change in bias (network reductions - subtractions)

# Chesapeake Bay Demo:

- Queried active, RT temperature gages (81)
- Want to examine ALL NLCD2016 land-uses
- But, want a “quick” number to show which variables are most/least representative!
  - “Bias”

## NHD Flowlines for Area of Interest

Location of Real-Time Temperature Gages 81 active

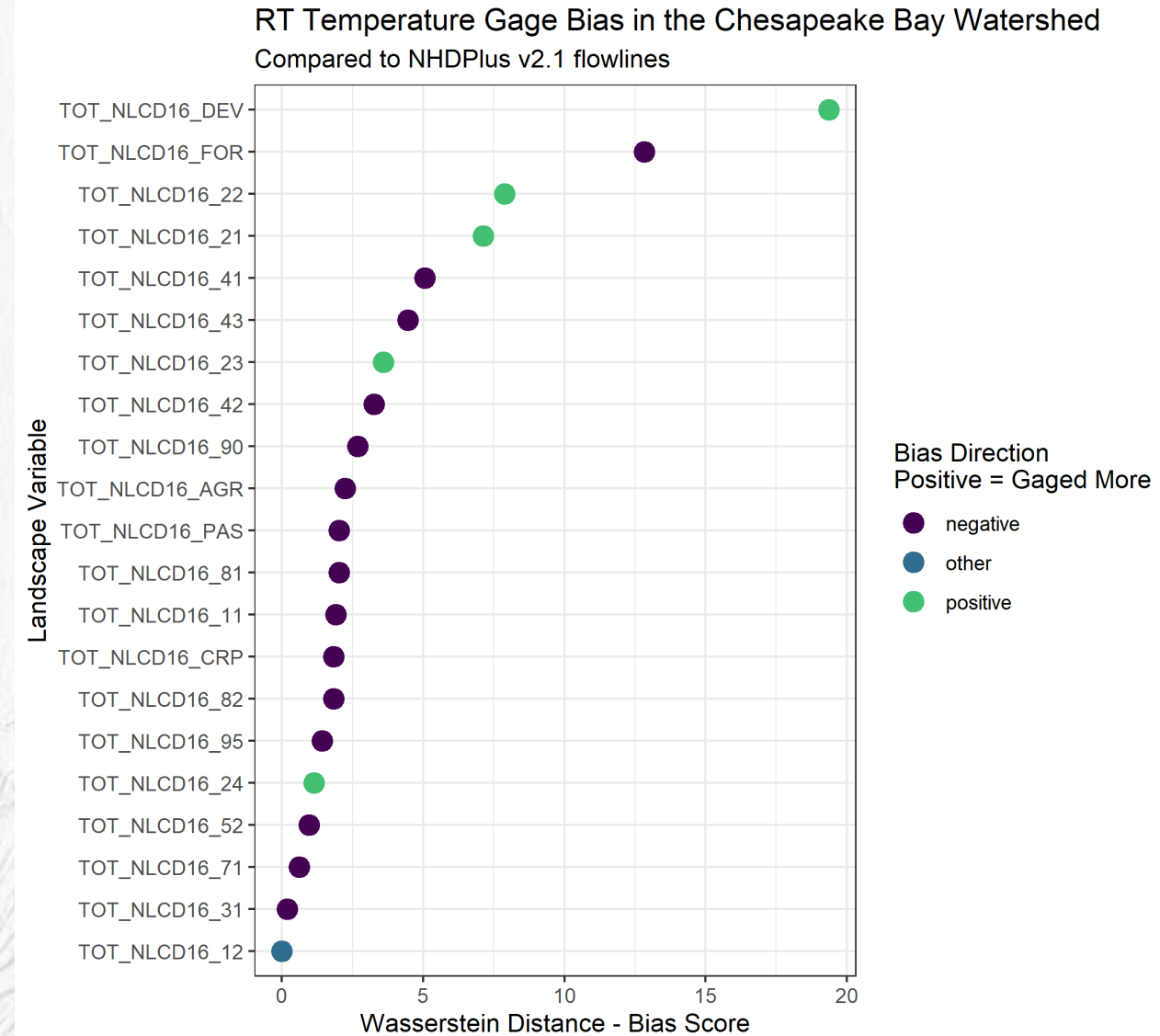
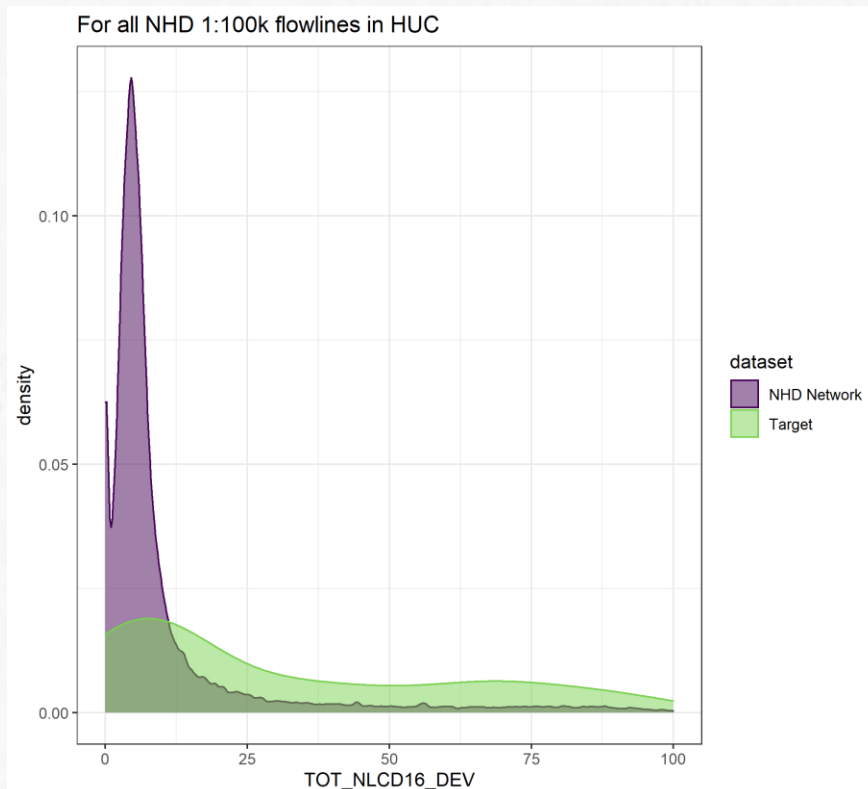




# Chesapeake Bay Demo: Bias Quantification

- Bias quantified using Wasserstein Distance ('earthmover score') for all NLCD2016 variables

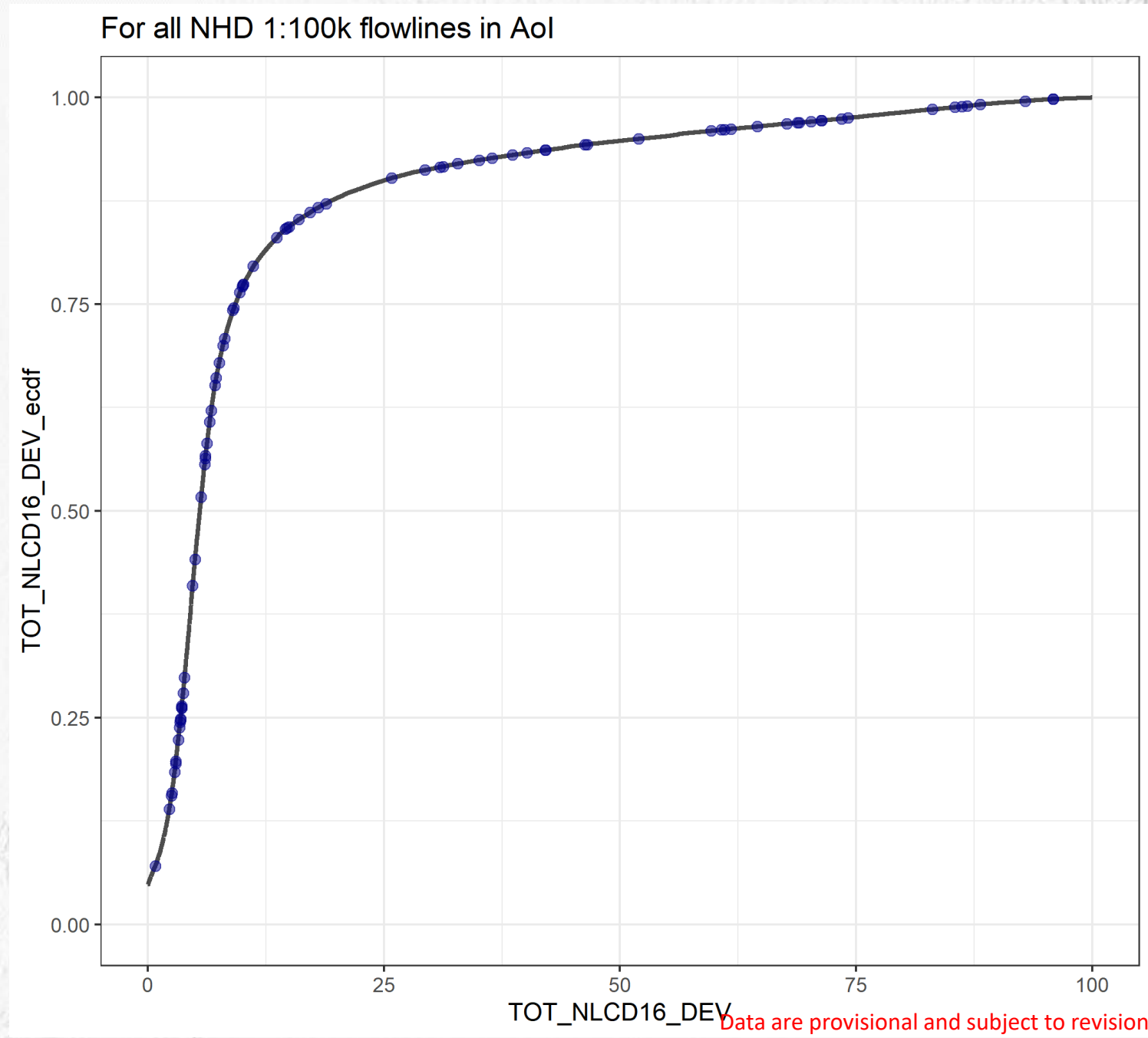
❖ “How much work is needed to “bulldoze” distributions to be the same?”



Data are provisional and subject to revision

# Gap Filling

- What if we're in the position to add gages?
- Distribution\_analysis()
  - Type: POINT



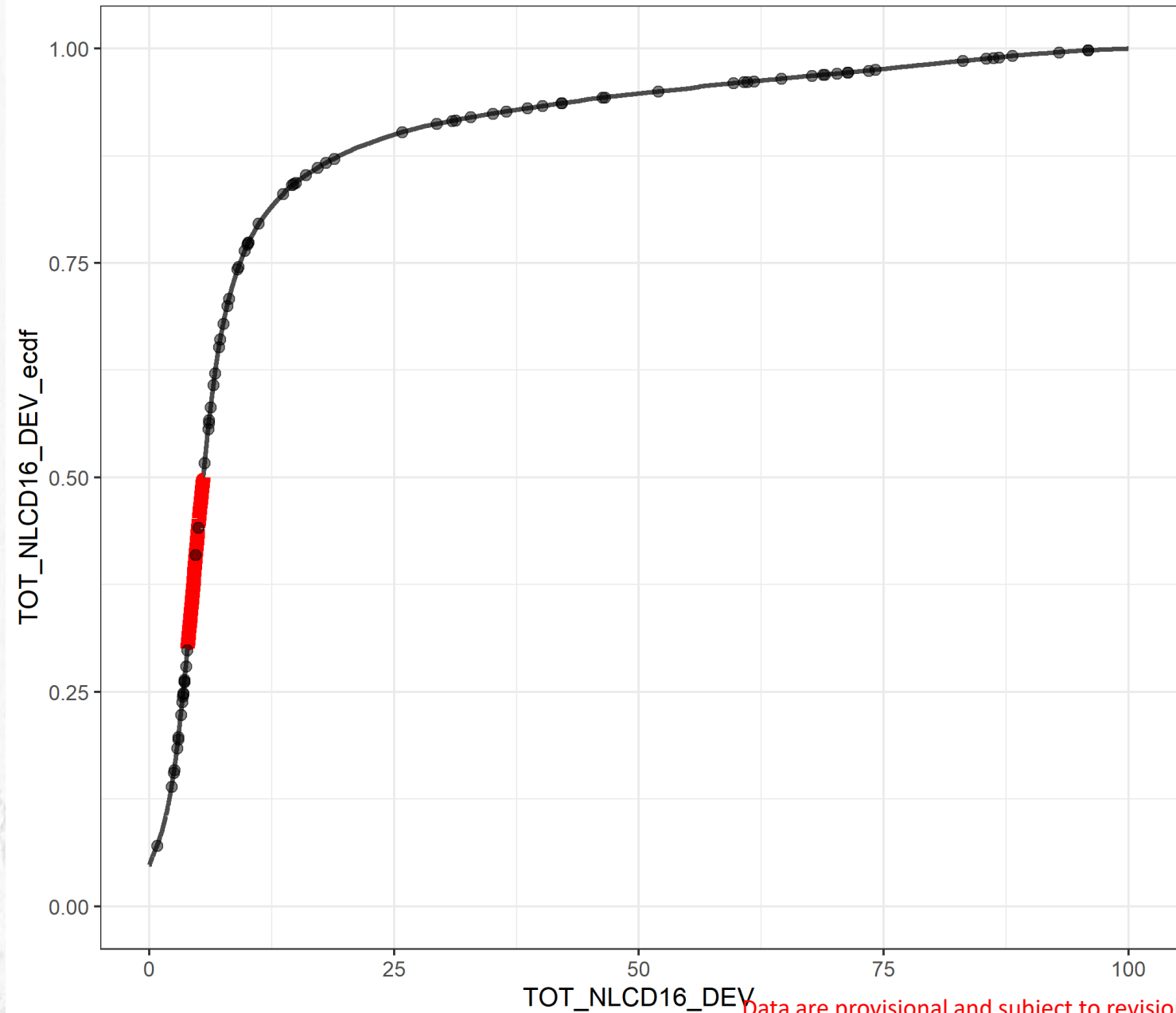
# Gap Filling

- Distribution\_analysis()
  - Type: POINT
    - Target GAP: 0.3-0.5

- First 100 Candidate COMIDs

- 4196310 4507182 4673153 4689163 4697207 4697237 4697673 4697675 4699853  
4710730 4764944 5890986 5892890 5893356 5893470 8101199 8101975 8102617  
8111937 8118019
- 8124969 8127489 8143572 8143584 8143586 8153211 8403479 8431706 8431922  
8431928 8436674 8456894 8457374 8477944 8478264 8501386 8507458 8539513  
8539919 8566657
- 8612147 9423769 10234695 2602267 4196050 4672055 4673205 4698827 4711598  
4763910 8100799 8110763 8112413 8112657 8124821 8126315 8126407 8139372  
8144200 8385881
- 8385899 8401523 8420616 8440195 8440839 8448752 8456514 8464839 8468791  
8490950 8493354 8493572 8501662 8503150 8526133 8539845 8552283 8574677  
8608481 8608793
- 9407856 9420343 14363248 14365368 8523575 8523579 8523595 8523749 8124983  
8143622 8151729 8457600 8457616 8479156 8479558 8547253 8609055 9423075  
9424555 22340163

For all NHD 1:100k flowlines in Aol



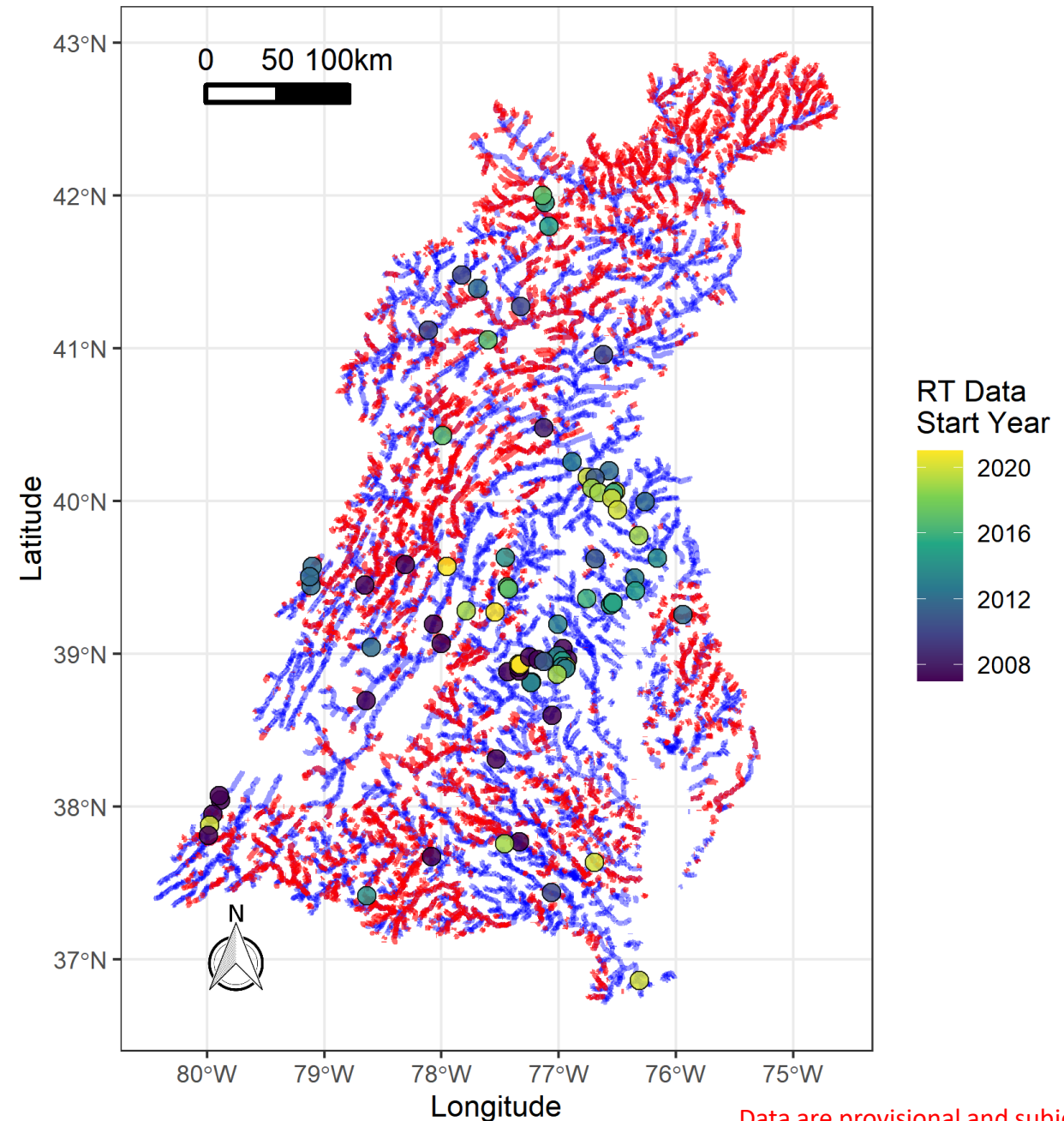
# Gap Filling

- Distribution\_analysis()
  - Type: POINT
  - Target GAP: 0.3-0.5

- First 100 Candidate COMIDs

- 4196310 4507182 4673153 4689163 4697207 4697237 4697673 4697675 4699853  
4710730 4764944 5890986 5892890 5893356 5893470 8101199 8101975 8102617  
8111937 8118019
- 8124969 8127489 8143572 8143584 8143586 8153211 8403479 8431706 8431922  
8431928 8436674 8456894 8457374 8477944 8478264 8501386 8507458 8539513  
8539919 8566657
- 8612147 9423769 10234695 2602267 4196050 4672055 4673205 4698827 4711598  
4763910 8100799 8110763 8112413 8112657 8124821 8126315 8126407 8139372  
8144200 8385881
- 8385899 8401523 8420616 8440195 8440839 8448752 8456514 8464839 8468791  
8490950 8493354 8493572 8501662 8503150 8526133 8539845 8552283 8574677  
8608481 8608793
- 9407856 9420343 14363248 14365368 8523575 8523579 8523595 8523749 8124983  
8143622 8151729 8457600 8457616 8479156 8479558 8547253 8609055 9423075  
9424555 22340163

NHD Flowlines for Area of Interest  
Location of Real-Time Temperature Gages



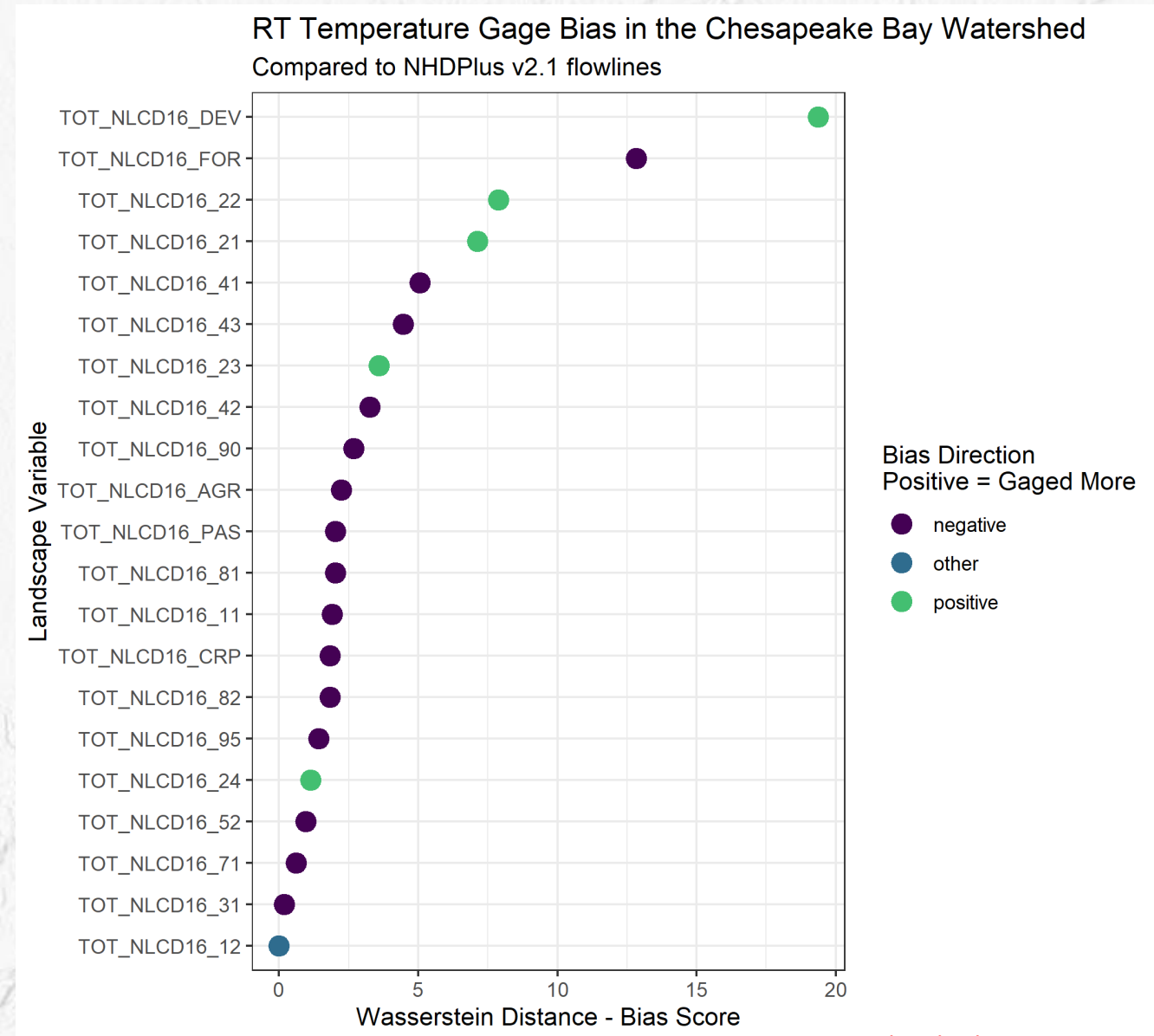
Data are provisional and subject to revision

# Gage Removal Analysis

- Minimizing the impacts of monitoring cuts
- Can bias scores evaluate impact of gage removal?

## Simulations:

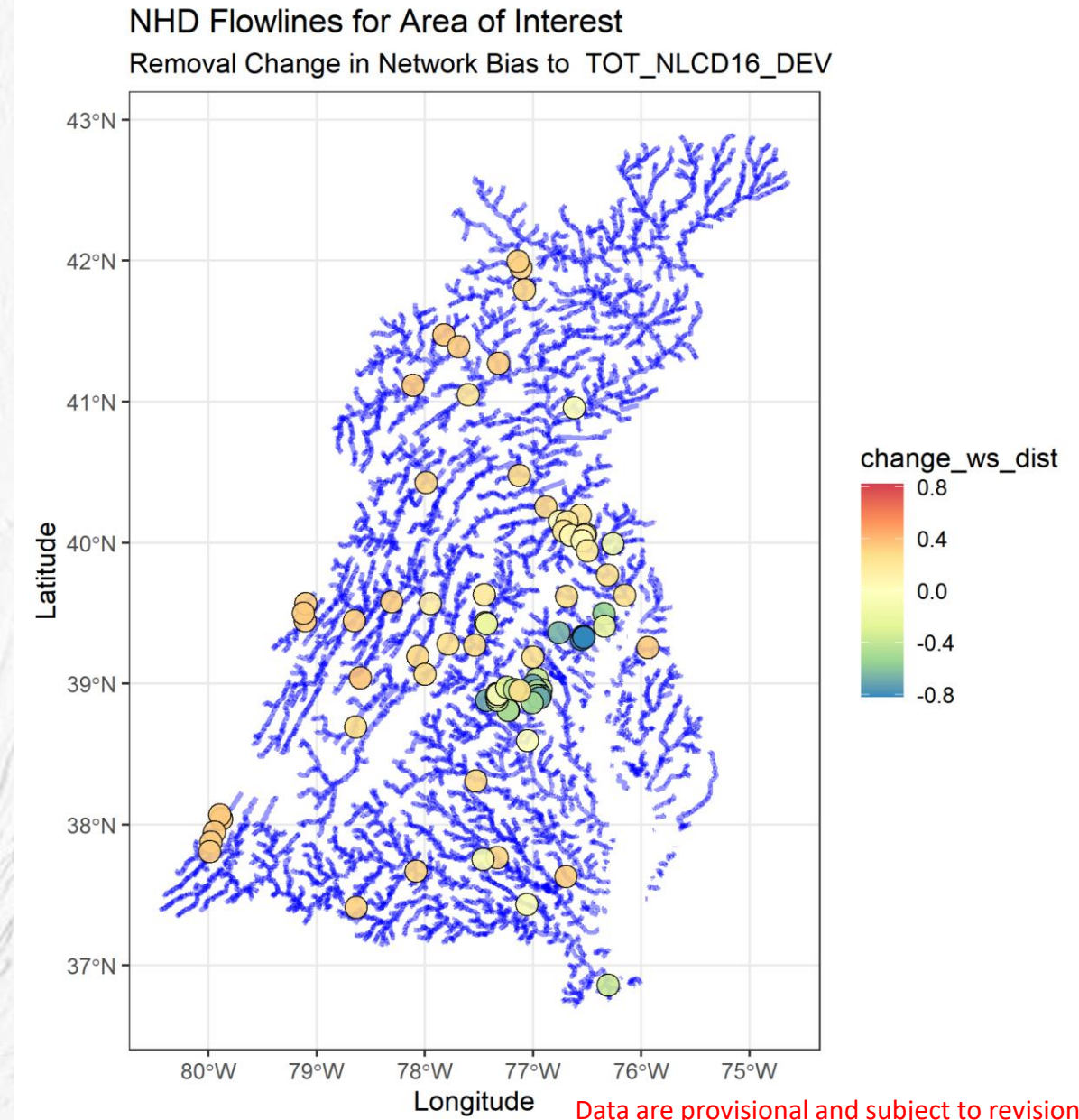
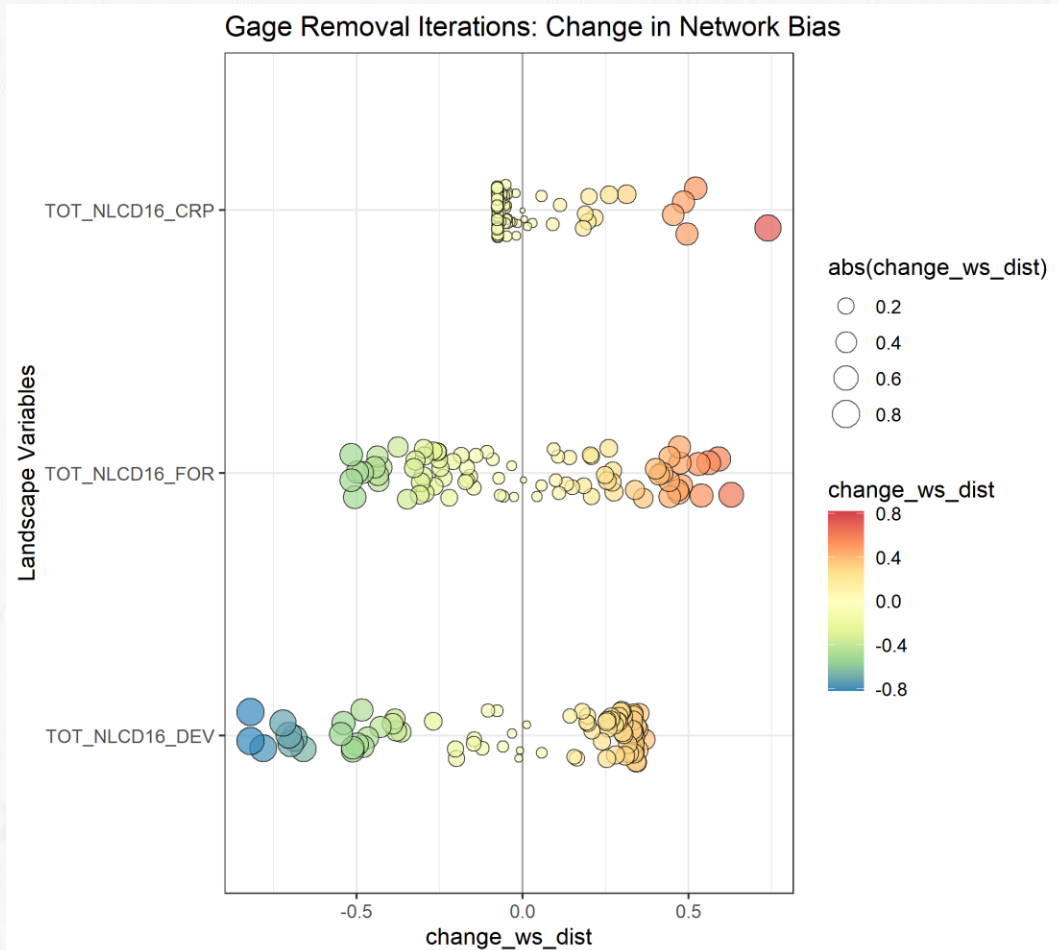
- Remove gage from network → recalculate bias
- Calculate change in bias from baseline
- Iterate for each gage in network



Data are provisional and subject to revision

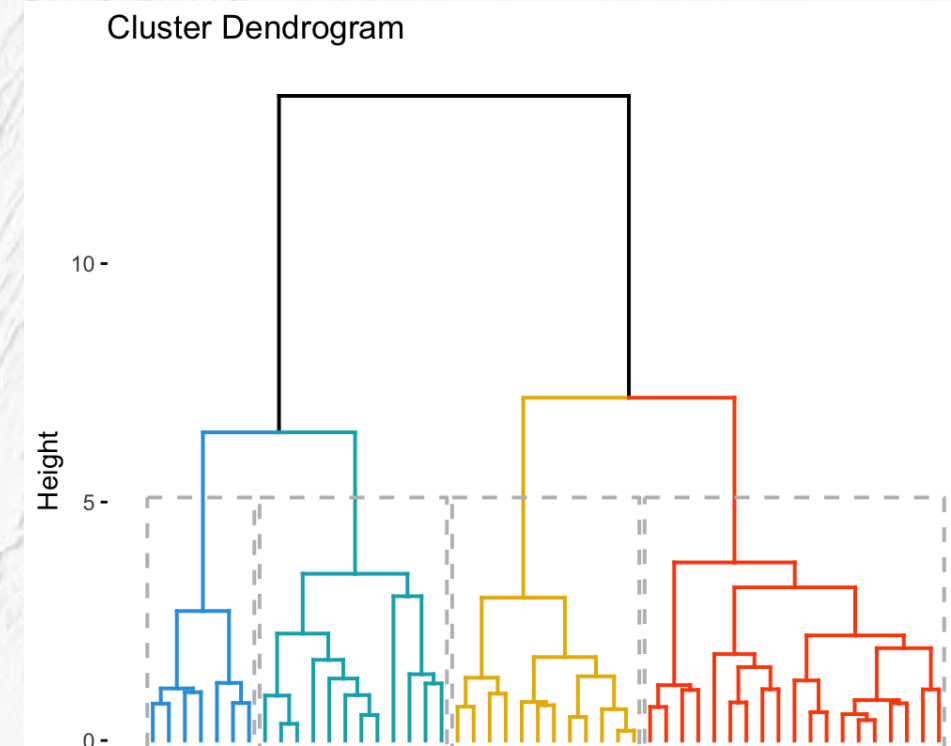
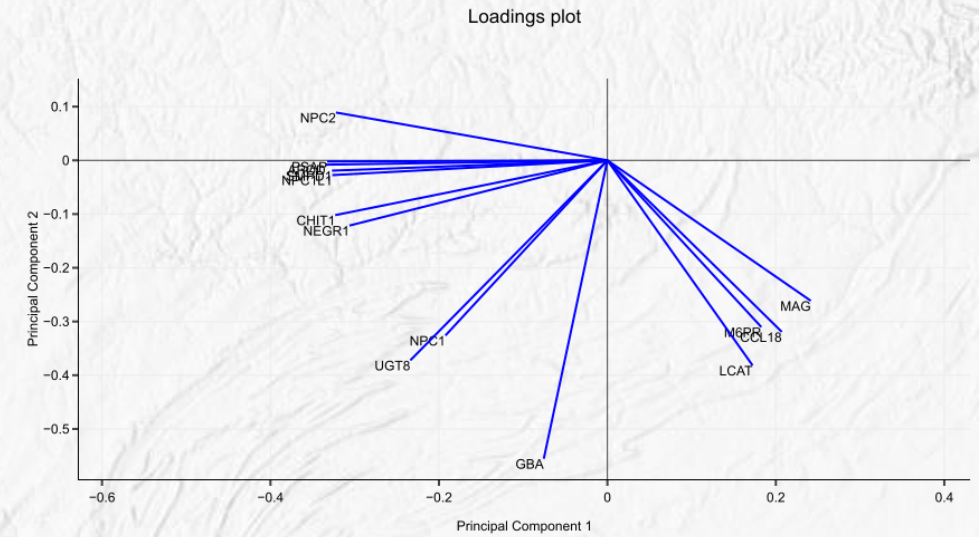
# Gage Removal Analysis

How does variable bias change if it was lost in the network?  
Results for every gage:



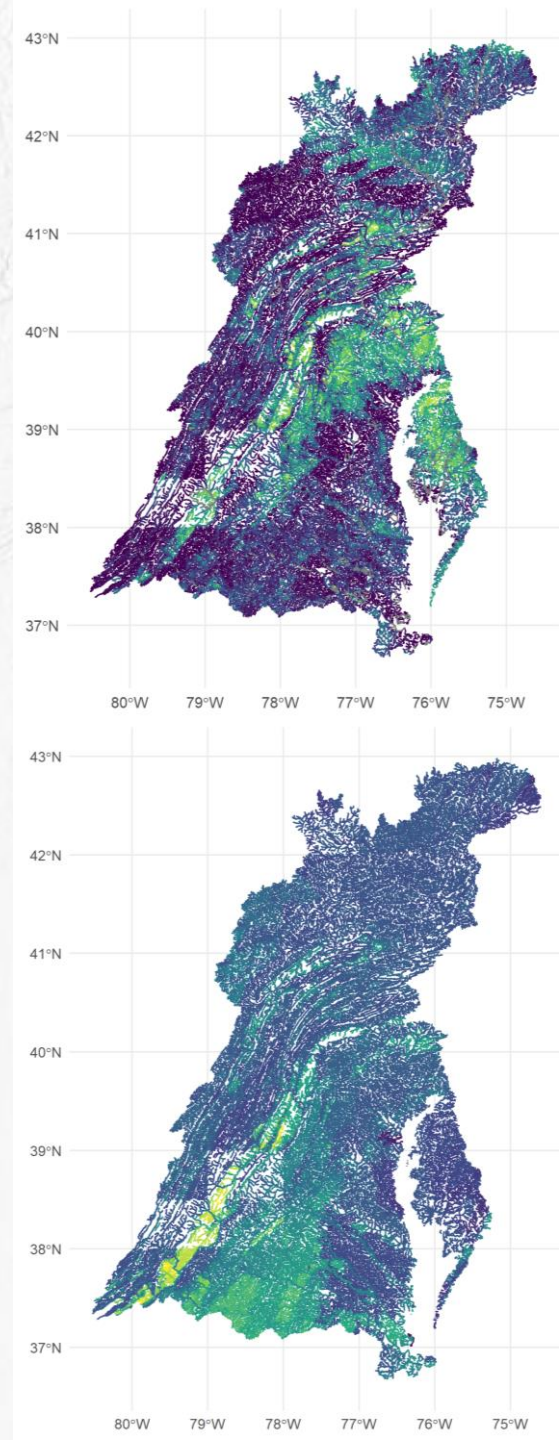
# Next steps

- Scripts are still in development
  - Bugs do exist (but it now works on mac)
  - Happy to have collaborators
  - New features continue to be added
    - Bias simulations with gage additions
- A way to evaluate multivariate bias
  - Equal weighting of considered variables?
  - Variable 'reduction' approach
    - PCA or hierarchical clustering?



# Tools are only tools

- These tools make network analyses “easier”
  - Handle “big” data
  - Avoid GIS pre-processing
  - Increase reproducibility
  - Built-in analysis methods
    - But always looking for more/improved methods!
- Yet cannot answer the fundamental issue:
  - Is the network representative?
    - *Representative of “what”?*





# What are your objectives and values?

- A network might be representative for one objective or goal, but not another
- With thousands of variables to evaluate, knowing what you are interested in and **why** is more important than ever
  - Decision-point is variable selection, not data availability
- Need a strong sense of priorities to evaluate inevitable tradeoffs
  - Network changes might decrease bias in one variable (or sensor type) but increase it in another



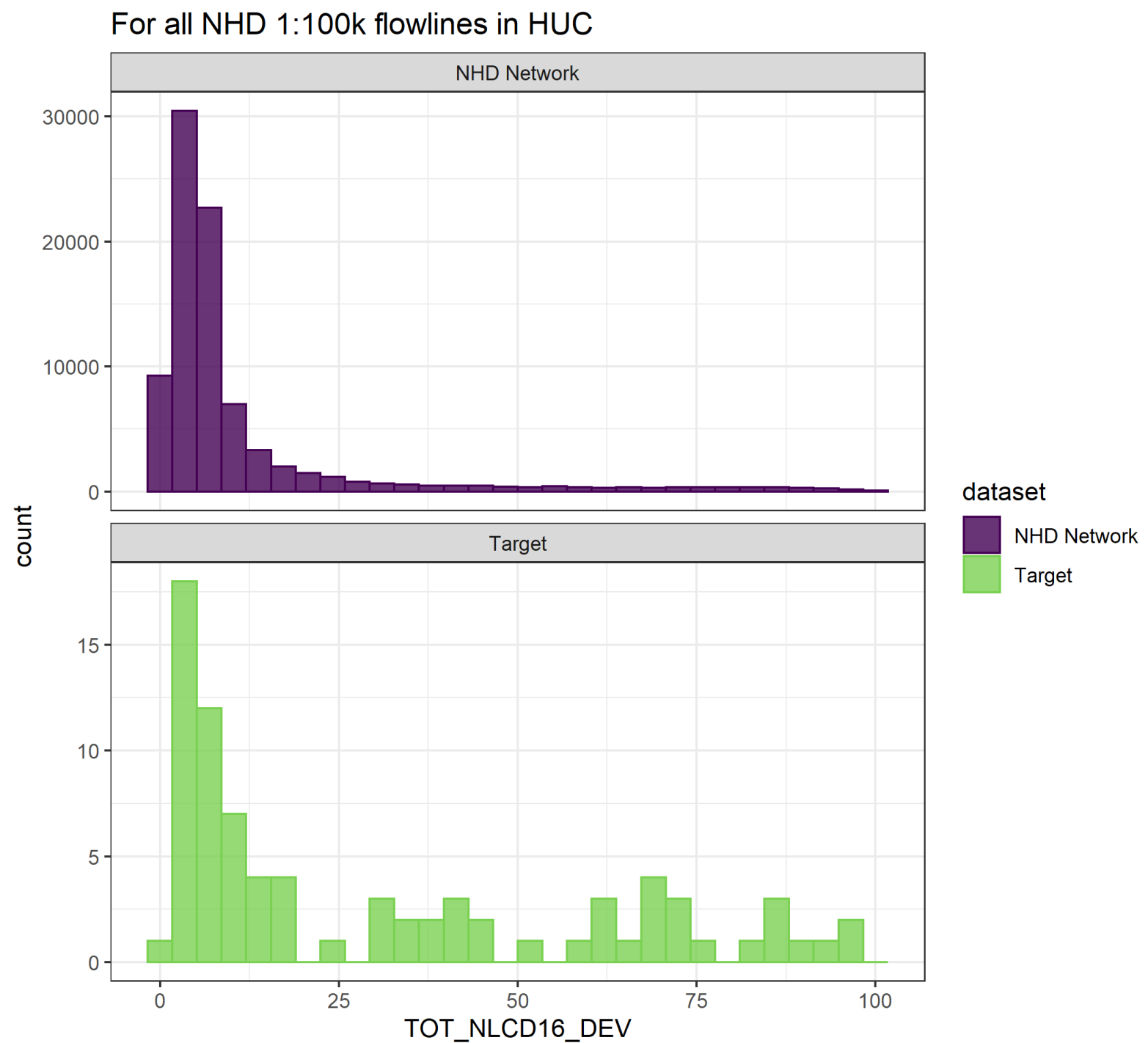
via Tok Talk

An aerial topographic map of a region, likely in the western United States, showing a network of rivers and a large reservoir on the right side. The terrain is depicted with green and brown tones, indicating elevation and vegetation. The text "Thank you!" is centered in a large, black, sans-serif font.

Thank you!

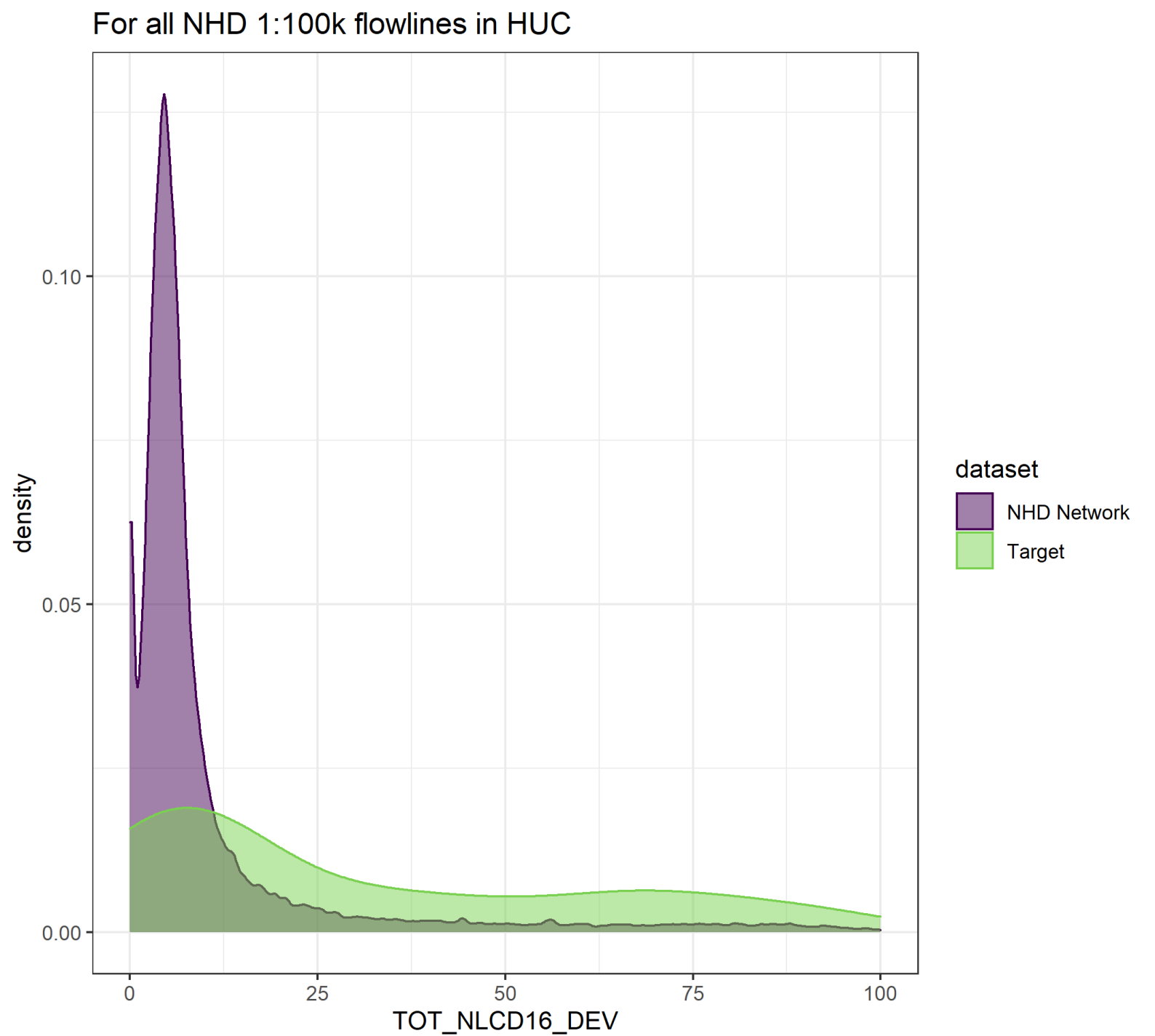
# Plotting

- Histogram\_analysis()



# Plotting

- Density\_analysis()



# Plotting

- `Distribution_analysis()`
  - Type: line

