# Considerations for Sharing, Using, and QC'ing Data

Cristina Mullin

Water Data Integration Branch

EPA Office of Water

# Learning Objectives

- To discover the potential data sources and formats that may be available from both the Tribe and other data partners for use in producing a water quality assessment

- To identify the factors that can affect the quality and suitability of data used for a water quality assessment

- To understand how to prepare data for analysis

# What Do We Want From Our Data?

- Inclusive: covering key parameters of concern
- Credible: to accurately reflect water quality conditions
- Robust: to reflect conditions under a variety of rainfall/flow regimes
- Useful: helping us identify appropriate solutions
- Efficient: the least cost for the most benefit!

# Considerations for Assessing Data

- Are there procedures for validating data?
  - Decision points to accept, reject, or qualify data
  - Procedures could include:
    - Examining results for high/low results
    - Checking calculations
    - Calculating precision & accuracy
      of instruments
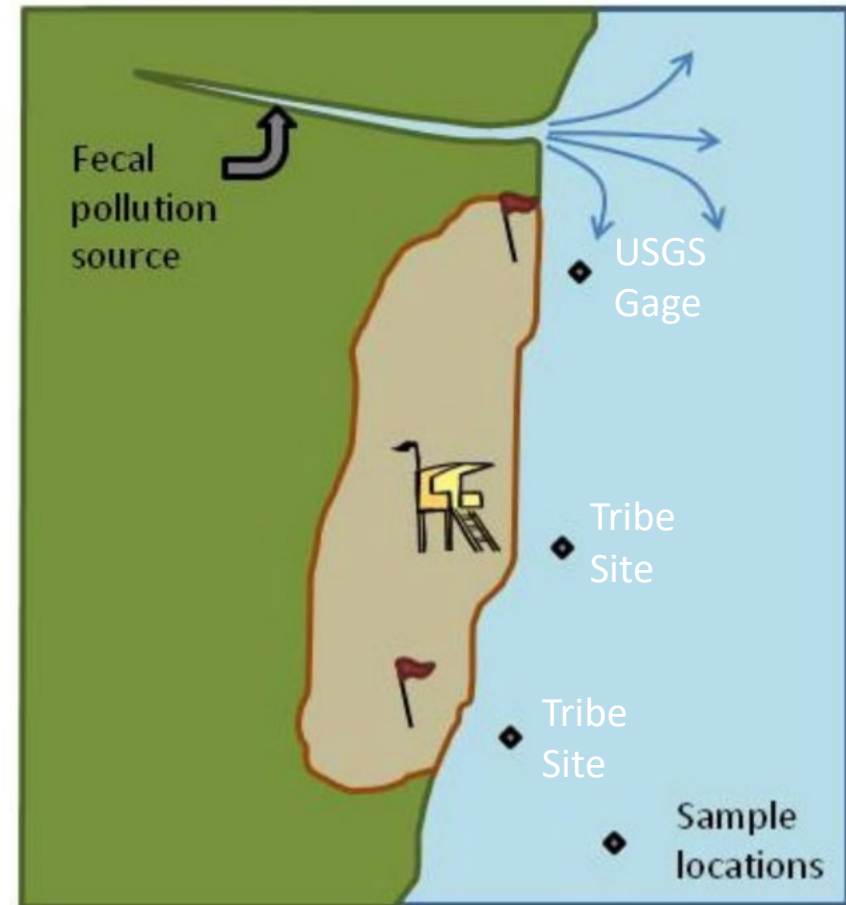- Are data adequate for a water quality assessment?

# A note about tribal data

- Tribal data collected with 106 funding must be shared with EPA at the end of each grant cycle (WQX/WQP).

- Tribal data collected using other resources does not have to be shared

# Why Consider Using Other Data?

- Might help to create a more comprehensive water quality assessment

- To fill data gaps

- To obtain other relevant information that supplements tribal data

- Important for tribes interested in TAS for Section 303(d)

- Supplement organizational monitoring for efficiency and cost savings
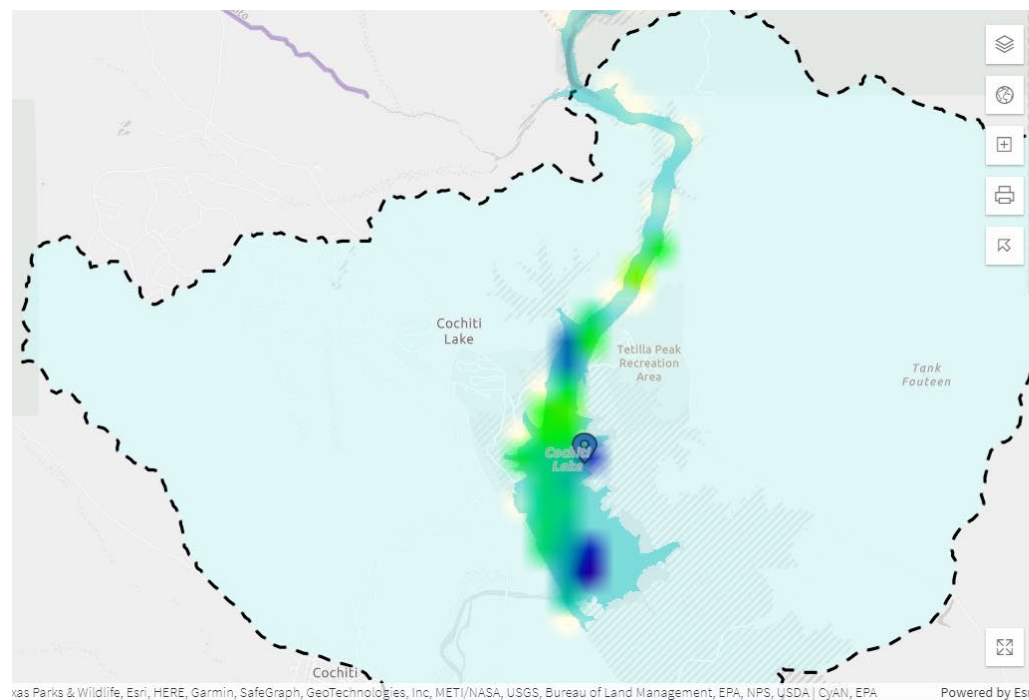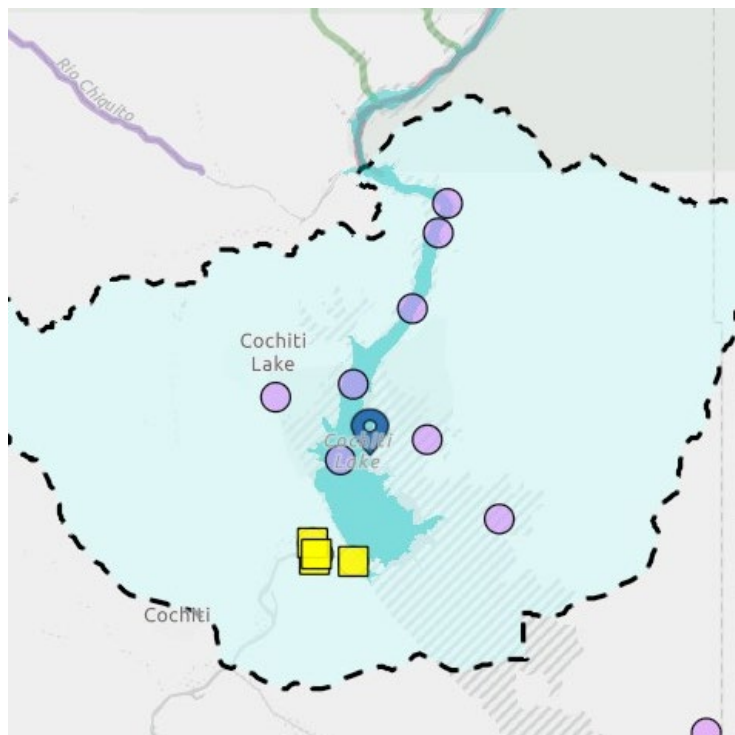
# What Other Types of Data Can Tribes Consider?

- Volunteer monitoring data
- Beach closure notices
- Fish consumption advisories
- Fish kills
- Source water assessments
- Waste site inventories
- Land use/cover data
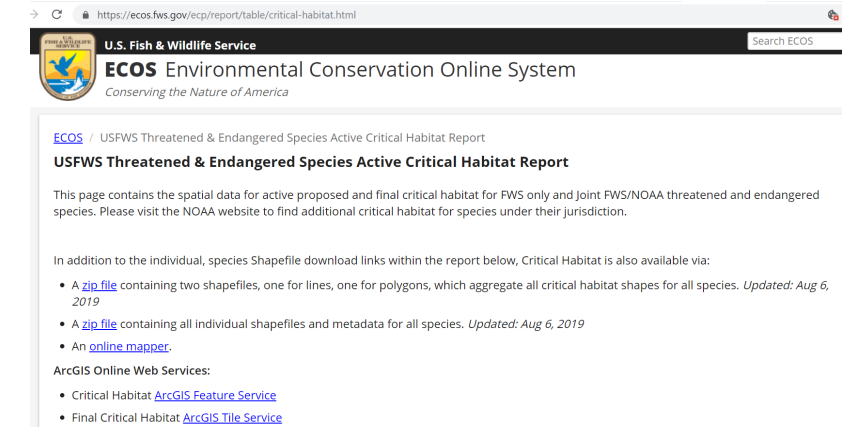- Hydrology, climate, geological studies/reports
- And more!

# Possible Sources for Additional Water Quality Data: Federal Agencies

- U.S. Environmental Protection Agency (ATTAINS, NARS)

- EPA & USGS Water Quality Portal (WQX/WQP)

- U.S. Geological Survey

- NASA Satellite Data

# Possible Sources for Additional Water Quality Data: Federal Agencies, Other Groups

- ## U.S. Fish and Wildlife Service
  - ### Fish, habitat

- ## U.S. Department of Agriculture Forest Service
  - ### Forest management plans

- ## National Oceanic Atmospheric Administration (coastal and estuarine data for both oceans and Great Lakes)

# Possible Sources for Additional Tribal Water Quality Data

- Bureau of Indian Affairs
- Indian Health Services
- Tribal commissions and ceded territory agencies
- Range of possible data
  - Water quality
    - Monitoring data
    - Fisheries (census and contaminant data)
    - Natural resources
    - Drinking water intake results
  - Source information
    - Septic systems
    - Landfills/waste sites

# Possible Sources for Additional Water Quality Data: State Agencies

- State Departments of environmental protection (305(b)/303(d) water quality assessment data, modeling, NPS assessments, source water protection assessments, watershed plans)

- Departments of natural resources (scenic rivers monitoring)

- Departments of health (recreational waters bacteria sampling, septic systems)



**State of New Mexico**
NONPOINT SOURCE
MANAGEMENT PROGRAM

2018 Annual Report

New Mexico Environment Department
Surface Water Quality Bureau
Watershed Protection Section

# Possible Sources for Additional Water Quality Data: Local Agencies

- Departments of Health
  - Septic system data
  - Beach monitoring data
- Water Utilities
  - Wastewater data
  - Drinking water monitoring data
- Soil and Water Conservation Districts
  - Water quality
  - Septic
  - Beach data

# Possible Sources for Additional Water Quality Data: Other Local Partners

- Universities
  - Research studies
  - Lab reports

- Watershed groups
  - Volunteer monitoring
  - Modeling
  - Nonpoint source assessments

# Online Data Sources/Tools

- Watershed Index Online: https://www.epa.gov/wsio

- Recovery Potential Screening: https://www.epa.gov/rps

- Healthy Watersheds Protection: https://www.epa.gov/hwp

- WATERS GeoViewer: https://www.epa.gov/waterdata/waters-geoviewer

- NOAA National Estuarine Reserve System-Wide Monitoring Program Data: https://coast.noaa.gov/digitalcoast/data/nerr.html

- How's My Waterway: https://mywaterway.epa.gov/

- National Aquatic Resource Surveys: https://www.epa.gov/national-aquatic-resource-surveys

# When Asking for Data From Potential Partners

- Be specific about data needs and intended use of the requested data

- Ask for the timeframe to fill data request

- Request a point of contact for follow-up

- Ask for the information needed to evaluate data quality
  - Metadata
  - QAPP

- Ask about the data format; be prepared to reformat

# Data Submission, Retrieval, and Assessment



**Water Quality Portal**

# Water Quality eXchange (WQX) Web and Nodes

- WQX Web
  - Sign up for account – https://cdx.epa.gov/
  - Requires less technical expertise and manual steps to upload
  - Allows you to submit data by uploading excel, .csv, or .txt files
  - Download Web Template Files and data rules
  - Find FAQs
- Custom submission application using WQX XML schema through Exchange Network Nodes or Node Clients
  - Allows you to submit data using coded files (.xml)
  - Custom WQX XML submission applications can be more efficient for organizations with larger databases and a need for automated submissions from internal databases
- WQX Domain Values  (data fields/metadata available)
- Group or 1-on-1 WQX trainings available (wqx@epa.gov)

# WQX hierarchy and terminology

- Organizations
  - All information is unique only to the organization
- User accounts
  - Multiple user accounts with different roles can be associated with an Organization
- Projects
  - Why you sampled
- Monitoring Locations
  - Where you sampled
- Activities and Results
  - Raw data referencing Orgs, Projects, & Monitoring Locations

# WQX QAQC Service

# Water Quality Portal (WQP)



**Water quality monitoring data is foundational to being able to answer important questions**

- Is my water safe?
- Is there enough?

**Format is the same for everyone who wants to share data**

- Water quality monitoring and data management is complicated
- Standardized, electronic data is more valuable than data in file cabinets (reusable, sharable, discoverable, interoperable, and includes important metadata)

**Usable data translates to knowledge, public awareness, and action**

- Reuse adds value!
- Supports CWA assessments and other water quality research
- Serves as the backbone for water data tools like HMW

# What data exists in the WQP?

WQX was modeled around the science, which does not change as much as IT software does, that's why it works

- Nutrients, metals, and biological data workgroups have made a lot of progress on best practices
- New WQX QAQC service for data submissions was deployed in spring 2022

**Characteristic/unit combos with the most data**

| CharacteristicName | Unit | RowCount |
|---|---|---|
| Temperature, water | deg C | 52859842 |
| Dissolved oxygen (DO) | ug/l | 11488857 |
| pH | None | 10306125 |
| Specific conductance | uS/cm | 8356335 |
| Salinity | ug/kg | 5055589 |
| Count | count | 4024377 |
| Conductivity | uS/cm | 3570617 |
| Turbidity | NTU | 3465425 |
| Temperature, air | deg C | 3436615 |
| Dissolved oxygen saturation | % | 2988042 |
| Inorganic nitrogen (nitrate and nitrite) | ug/l | 2868888 |
| Total suspended solids | ug/l | 2843607 |
| Phosphorus | ug/l | 2551290 |
| Depth, Secchi disk depth | in | 2027629 |
| Depth | in | 2002344 |
| Kjeldahl nitrogen | ug/l | 1987872 |
| Orthophosphate | ug/l | 1823127 |
| Count | % | 1683594 |
| Organic carbon | ug/l | 1655437 |
| Chloride | ug/l | 1549265 |

8/9/21

# Retrieving WQP data from the WQP

WQP Web Interface:
https://www.waterqualitydata.us/

- WQP Demo on How to Download Data (2015)

- WQP Demo on How to Download Data (2019)

How's My Waterway

TADA

# Screening/Flagging Data

- Determine data credibility
- Ensure appropriate data quality
- Potential screening criteria
  - Parameters and methods used
  - Range checks

# Data cleaning - breadth

- Join multiple WQP profiles if needed
- Review metadata & Filter
  - Filter by media type and waterbody type, etc.
- Unit conversions
- Synonym checks
- Duplicates

Which of these we need to do depends on the dataset...

# Data cleaning - breadth

- Correct data types (numeric or categorical), address symbols in results
- Speciation considerations
- Check if monitoring equipment /methods changed over time
- Quality checks
  - Outlier detection
  - Location accuracy

Which of these we need to do depends on the dataset…

# Data cleaning - breadth

- Staff changes may impact the way metadata is reported (sometimes things get into different WQX fields by accident)

- Different organizations may list data differently in WQP even though WQX tries to prevent that

Which of these we need to do depends on the dataset...



Drop unwanted columns

Handle missing data

Correct data types

Merge multiple datasets

Data Cleaning

Drop the duplicates

Review metadata and filter dataset

Convert units

# Non-detections

- Set Non-Detections equal to the Limit of Detection
- Set Non-Detections equal to the 1/2 times the Limit of Detection
- Set = 1/x detection limit (you define x)
- Assign values to Non-Detections using the Kaplan-Meier method

# Interquartile Range (IQR)



- Method to identify data that are different than approximately 99% of the data available for the assessed parameter
  - Upper Outlier = 75th Percentile + 1.5 * (75th percentile - 25th percentile)
  - Lower Outlier = 25th Percentile - 1.5 * (75th percentile - 25th percentile)

- May want to flag data that falls above or below the upper or lower value

# Data Screening Considerations: Parameter v. Methods

- Parameters can have many forms (total nitrogen, total Kjeldahl nitrogen, nitrate, nitrite )

- Essential to specify the chemical form of the parameter

- An analytical method is the procedure for determining the amount/concentration of the parameter
  - Several analytical methods can apply to a parameter
  - Essential to specify which analytical method is used
  - Limits of detection are also important to consider. Specifically, when the water quality standard is near the detection limit.

# Data Screening Considerations: Range Checks

- Identify the range of possible concentrations for each parameter based on:
    - Site
    - Historical data
    - Recent watershed changes
- Values outside of that range may be in error
    - Investigate upstream/upland conditions before discarding data
    - Check to see if the collection method requires field blanks and make sure they are all below the limit of detection (indicates whether sample is contaminated or not)

# Data Screening Considerations: Box plot

# Data Screening Considerations: Histogram



How is the data distributed?

# Total Phosphorus_as P_ug/L

# Data Screening Considerations: Scatter Plots



Scatter Diagram

# Time series – type of Scatter plot

# Depth profile - type of Scatter plot

# Exercise: Data Screening

# Which datapoints need further review?

| Date | pH (standard units) | Field comment |
|------|------|------|
| June 1 | 6.9 | Cloudy |
| June 14 | 7.1 | |
| June 23 | 6.8 | Sunny but cool |
| July 8 | 5.2 | |
| July 15 | 7.1 | Windy |
| July 20 | 7.1 | |
| July 29 | 7.0 | Overcast |
| August 2 | 6.9 | |
| August 8 | 6.8 | No pH 7 calibration solution |
| August 16 | 7.1 | |
| August 23 | 8.2 | Drizzling |
| August 31 | 7.2 | |

Which datapoints still need further review?

| Date | pH (standard units) | Field comment |
|---|---|---|
| June 1 | 6.9 | Cloudy |
| June 14 | 7.1 | |
| June 23 | 6.8 | Sunny but cool |
| July 8 | 5.2 | |
| July 15 | 7.1 | Windy |
| July 20 | 7.1 | |
| July 29 | 7.0 | Overcast |
| August 2 | 6.9 | |
| August 8 | 6.8 | No pH 7 calibration solution |
| August 16 | 7.1 | |
| August 23 | 8.2 | Drizzling |
| August 31 | 7.2 | |

| Date | pH (standard units) | Field comment |
|---|---|---|
| July 7 | 7.0 | Began to rain after sampling |
| July 9 | 6.9 | |
| July 11 | 7.0 | |
| July 13 | 7.1 | Windy |
| July 15 | 7.1 | Windy |
| July 17 | 7.0 | |
| July 19 | 7.1 | |
| August 3 | 6.9 | |
| August 8 | 6.8 | |
| August 13 | 6.9 | |
| August 18 | 7.1 | |
| August 24 | 8.0 | Limestone gravel in pile near stream |
| August 28 | 7.6 | |

# Which datapoints need further review?

| Date | Total phosphorus (mg/L) | Total nitrogen (mg/L) |
|---|---|---|
| June 15, 2012 | 1.9 | 6.7 |
| July 16, 2012 | 1.7 | 6.4 |
| August 13, 2012 | 2.1 | 8.2 |
| September 14, 2012 | 2.3 | 8.1 |
| May 4, 2017 | 0.751 | 4.53 |
| May 24, 2017 | 0.813 | 4.44 |
| June 15, 2017 | 0.795 | 4.83 |
| July 5, 2017 | 0.702 | 4.61 |
| July 26, 2017 | 0.699 | 4.45 |
| August 14, 2017 | 0.785 | 4.56 |
| August 31, 2017 | 0.803 | 0.43 |
| September 19, 2017 | 0.797 | 4.42 |
| May 1, 2018 | 0.789 | 3.34 |
| May 16, 2018 | 0.812 | 3.42 |
| June 2, 2018 | 7.78 | 3.53 |
| June 18, 2018 | 0.808 | 3.67 |
| July 2, 2018 | 0.825 | 3.79 |
| July 15, 2018 | 0.837 | 3.77 |
| July 29, 2018 | 0.914 | 3.45 |
| August 10, 2018 | 0.956 | 3.51 |
| August 26, 2018 | 1.002 | 3.62 |
| September 10, 2018 | 3.6 | 0.998 |
| September 23, 2018 | 0.923 | 3.54 |

# Which datapoints need further review?

| Date | Turbidity (NTU) | Field Comments |
|------|-----------------|----------------|
| July 10, 2019 | 10 | Slight drizzle |
| July 25, 2019 | 12 | Light rain |
| August 11, 2019 | 9 | Clear |
| August 21, 2019 | 29 | Rain in the morning; bankfull flow |
| September 8, 2019 | 10 | None |
| September 19, 2019 | 11 | None |
| October 4, 2019 | 16 | Cattle in field near stream |
| October 18, 2019 | 22 | Banks appear to be trampled |
| November 1, 2019 | 26 | None |
| November 15, 2019 | 28 | Cattle in field near stream |
| November 31, 2019 | 26 | None |
| December 4, 2019 | 23 | None |

# Which datapoints need further review?

| Site X-42 Sample Date | Lab E. Coli | Lab Comments |
| --- | --- | --- |
| July 10, 2019 | 280 | No field notes |
| July 25, 2019 | 210 | None |
| August 11, 2019 | 160 | Started using new sample bottles |
| August 21, 2019 | 190 | Chain of custody form not signed |
| September 8, 2019 | 240 | New sampling staff |
| September 19, 2019 | 760 | Holding time exceeded by 3 hours |
| October 4, 2019 | 250 | None |
| October 18, 2019 | 180 | Duplicate |
| October 18, 2019 | 190 | Duplicate |
| November 1, 2019 | 210 | None |
| November 15, 2019 | 690 | Sample not on ice |
| November 31, 2019 | 200 | None |
| December 4, 2019 | 190 | None |

# Organizing Your Data for Analyses

- Entering your data on a spreadsheet GREATLY simplifies the analysis
  - It also helps to protect and preserve your data
- Clean up your data by:
  - Making sure everything is consistent, such as dates, parameter names, site designations, etc.
  - Checking for commas vs. decimal points
  - Looking for letters within numbers
- Keep data organized via:
  - Filename protocols with dates, controls on data entry, periodic reviews

# Bottom Line in Assessing Data Quality

- Identify the data being considered for use
  - Tribal (primary)
  - Non-tribal  (secondary)
- Collect information on how the data was produced (sample collection, analysis, reporting procedures)
- Review data quality guidance used in producing the data (QAPP, DQOs/DQIs)
- Screen the data for obvious problems
  - Poor documentation of procedures
  - Values below detection limits, significant outliers, etc.
- Evaluate the usefulness of the data
- Document justifications for data use / non-use

# Data Quality Scenario 1:

- Watershed group collects biweekly chemistry samples
  - June through  September
  - Purpose: evaluation of effects on macroinvertebrate health during summer low-flow critical conditions
- An upper Midwest Tribe wants to use data set to estimate annual pollutant loading
- Discussion: Is this watershed data representative of the conditions the tribe wants to evaluate for their water quality assessment? Why or why not?
- *(HINT: When might pollutant loading be highest and when is it lowest and what data did you capture?)*

# Data Quality Scenario 2

- Watershed group worked with trained volunteers to collect water quality data
  - Used field test kits
  - Purpose: To determine the concentration of a specific pollutant to the nearest milligram per liter
- Tribe's data is analyzed in a lab to the nearest microgram/liter
- Discussion: How might the tribe use both datasets for the water quality assessment? What additional information might be needed?

# Key Take-Aways

- Identify all existing and readily available data for the assessment
    - Parameters collected by the tribe through its monitoring program
    - Other relevant data and information about tribal waters or watershed
- Use online data tools and work with other local data partners
- QAPPs and DQOs are foundational to assessing data quality
- Evaluate all existing and readily available data for the water quality assessment
    - Review for quality through QAPP review and data screening
    - Use only data of adequate quality after review and screening