# Utah's Use of Cloud-Based Big Data Tools for Continuous Water Quality Data

Paul Burnett and Alan Ochoa
Utah Division of Water Quality

UTAH DEPARTMENT *of*
ENVIRONMENTAL QUALITY
**WATER QUALITY**

# Utah Division of Water Quality Continuous Water Quality DB Team

Alex Heppner

Ben Holcomb

Alan Ochoa

Ryan Parker

Marshall Baillie

Julian Carroll

Toby Hooker

Paul Burnett

# Project Overview

- Goal: Centralize and manage continuous water quality data from multiple monitoring sites and sources throughout Utah.
- Scope:
  - Multiple sensor types (e.g. lake, stream, piezometers)
  - Multiple parameters (e.g. water temp, ph, etc)
  - Real-time or near-real-time data streams
  - Bulk data ingestion from deployed loggers
  - Allow multi-platform access for analysis
- Users: Program staff, researchers, and partner agencies, general public.

# Constraints

Out-of-the-box solutions are too expensive

No dedicated staff for database development or management

Existing system pulling data from streaming buoys was manual (monthly)

No existing framework for logger data (~1200 datasets and counting)

# System Architecture

**Data streams**
- Deployed standalone loggers
- Streaming data from buoys
- Streaming USGS data

**Leveraging Google Sheets**
- Location Table
- Samping Event Table
- Bulk Data Entry (template)
- Google forms for basic data entry
- Apps Scripts for automation

- Bigquery for data warehousing
- Cloud Functions for automation
- Scheduled queries for:
  - Automated data wrangling
  - Denormalization for visualization
- Looker Studio for visualization
- Google Shared Drive - Working Folders for Raw Data

| **Location** |
| --- |
| Location_ID |
| Geo Location |
| Water Name |
| Location Type |
| Other info (HUC12, AU), etc |

| **Sampling_Event** |
| --- |
| Sample ID |
| Location ID |
| Collecting Org |
| Logger Info |
| Added to BQ? |

| **Water_Data** |
| --- |
| Sample ID |
| Meas Type |
| Units |
| Measurement |
| Flags |

## HighFrequencyLocation

| Column | Type | Description |
|---|---|---|
| HF_Location_ID | Integer | Unique ID of locations |
| Logger_Location_Description | Text | Description of Location |
| Location_Notes | Text | Additional Location Notes |
| Lat | Number | Decimal Degrees |
| Long | Number | Decimal Degrees |
| USGS_HUC12 | Text | 12-Digit HUC |
| State | Text | State |
| County | Text | County |
| Water_Name | Text | Water Name |
| MLID | Integer | Cross reference to Utah monitoring locations |
| Date_Added | Date | Date that the location record was added. This is just for record keeping. |
| Add_Log | Text | The name of the person who added the record. |
| Site_Type | Text | Stream, lake, reservoir, canal |
| Folder_Link | Link | Each location has a google folder containing subfolders of raw data. |

Notes:
- Location, Sampling Event, qc_log and other "small tables" managed in google sheets. The sheets are linked within bigquery as connected, external tables.
- Location table links locations and attributes to a geo point.
  - Includes an MLID crosswalk column to link to AWQMS.
- Sampling Event table tracks info about each deployment.
- Auto incrementing and categorical data values are controlled by using google forms for data entry.
- Folder structure for raw data folders follows the Location -> Sampling Event structure.
- All actual water data is stored in the Water_Data table.
  - Separate water data tables store streamed data and datasets with older qaqc_Level values.
  - Primary identifyer is Sampling Event ID.
  - Key value is a combination of Sampling_Event, DateTime_UTC, Measurement_Type, qaqc_Level
- DateTime_Local is the local time stored consitently as UTC minus 6 hours - Mountain Daylight Time.
- QC is applied to a Meas Type for entire sampling event.
  - Only increase qc level if measurement values are changed.
  - Adding qc flag would not constitute a new qc level.

## HighFrequencySamplingEvent

| Column | Type | Description |
|---|---|---|
| Sample_ID | Integer | Unique ID of each sampling event |
| Location_ID | Integer | Links the sampling event to each location |
| Sample_Year | Integer | Four-digit year of when the device was deployed |
| Logger_ID | Text | Serial number of the deployed device |
| Logger_Type | Text | Type of device such as standalone logger, telemetered device, or USGS gage |
| Collecting_Org | Text | The organization that deployed the device |
| Comment | Text | Comments about the device deployment |
| Date_Added | Date | Date that the sampling event record was added. Used for record keeping |
| Add_Log | Text | The name of the person who added the record |
| Raw_Data_Folder_Link | Link | The google folder for each sampling event. These are subfolders under each location |
| createNewFolder | Text | Used as a trigger for an apps script to create a new folder |
| wqDateLiveID | Integer | Links to the device id for buoys deployed under wqDataLive |
| BigQuery_Entered | Text | Tracks the progress of data entry into the BigQuery database |

## Water_Data

| Column | Type | Description |
|---|---|---|
| Sampling_Event | Integer | Links to the sampling event ID in the Sampling Event table. |
| Datetime_UTC | Timestamp | Timestamp |
| Datetime_Local | DateTime | UTC minus 6 hours (i.e. MDT) |
| Water_Measurement | Float | Measurement Value |
| Unit_of_Measurement | Integer | Coded value of the unit of measurement |
| qaqc_level | Integer | This column tracks the version of the data. By default all data coming in to the database has the value of 1. This value increases if qc adjustments are applied to the measurement type at the dataset level. |
| Measurement_Type | Integer | Coded value of the water quality parameter being measured. Links to values in a lookup table. |
| Meas_Flag | Integer | Used to flag individual values as erroneous. Links to values in a lookup table. |

## tblQC_log

| Column | Type | Description |
|---|---|---|
| Sample_Event | Integer | Links to the sampling event of the dataset. |
| Measurement_Type | qcLevel | The parameter of the |
| qaqc_level_original | Integer | The original qc level of the data modified |
| qaqc_level_updated | Integer | The updated qc level of the modified data |
| qc_Date | Date | This is the date that the QC m |
| qc_Narrative | Text | Describes the qc adjustements made. |

# Bigquery for Big Tables



## Water Data Tables

- Data from dataloggers and streamed data
- QC archived data

## USGS Data tables

- Final data table (data older than 4 months)
- 4 month data (provisional data)
- Sites and parameter codes
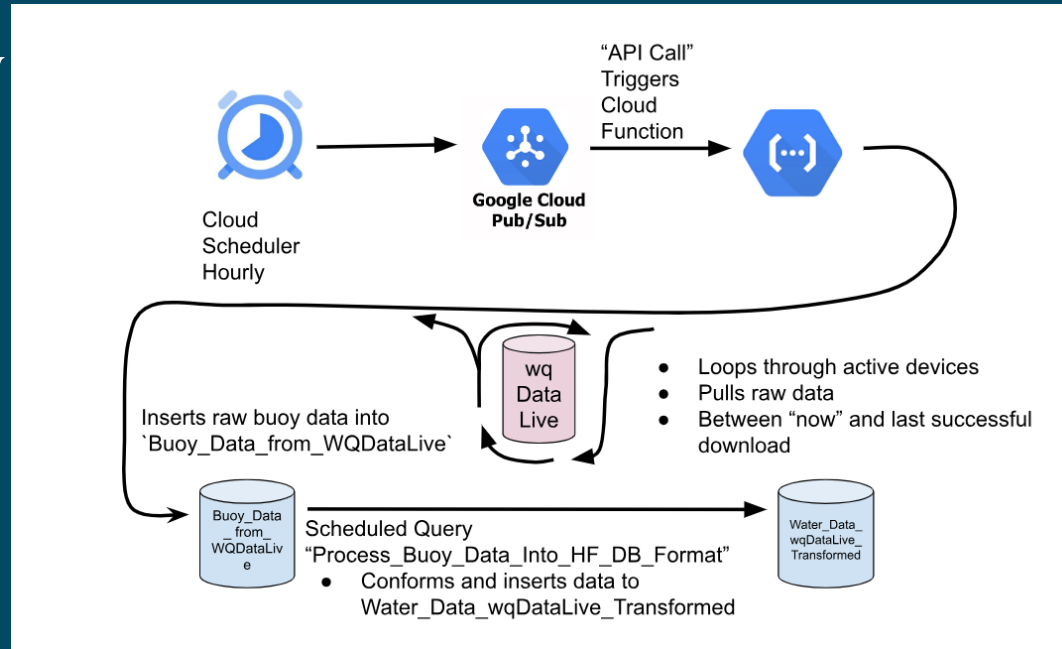
## External sheets

- linked to google sheets

## Flat Views

- Data flattened for dashboards
- Pulled from data tables hourly/daily with scheduled queries

# Cloud Functions for Automation

Lightweight, event-driven functions that run backend code in response to events without needing to manage servers.

- Using for streaming buoy data from wqDataLive

- Streaming USGS data

- Testing of streaming microcontroller data (soon)

# Public Data Accessibility
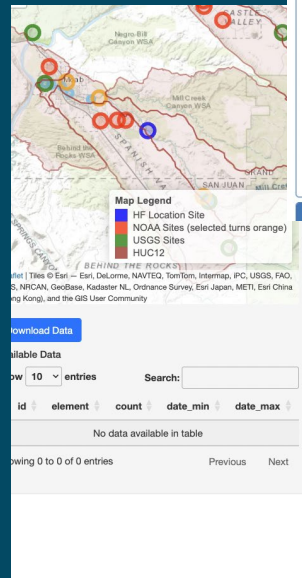
# Seamless Access to Continuous Data

- Pull millions of rows with just a few lines of R code:

```
33
34   query = "SELECT * FROM `exampleProject.dataset.exampleLocations`"
35   location_types <- bqQuery("exampleProject",query)
36
```
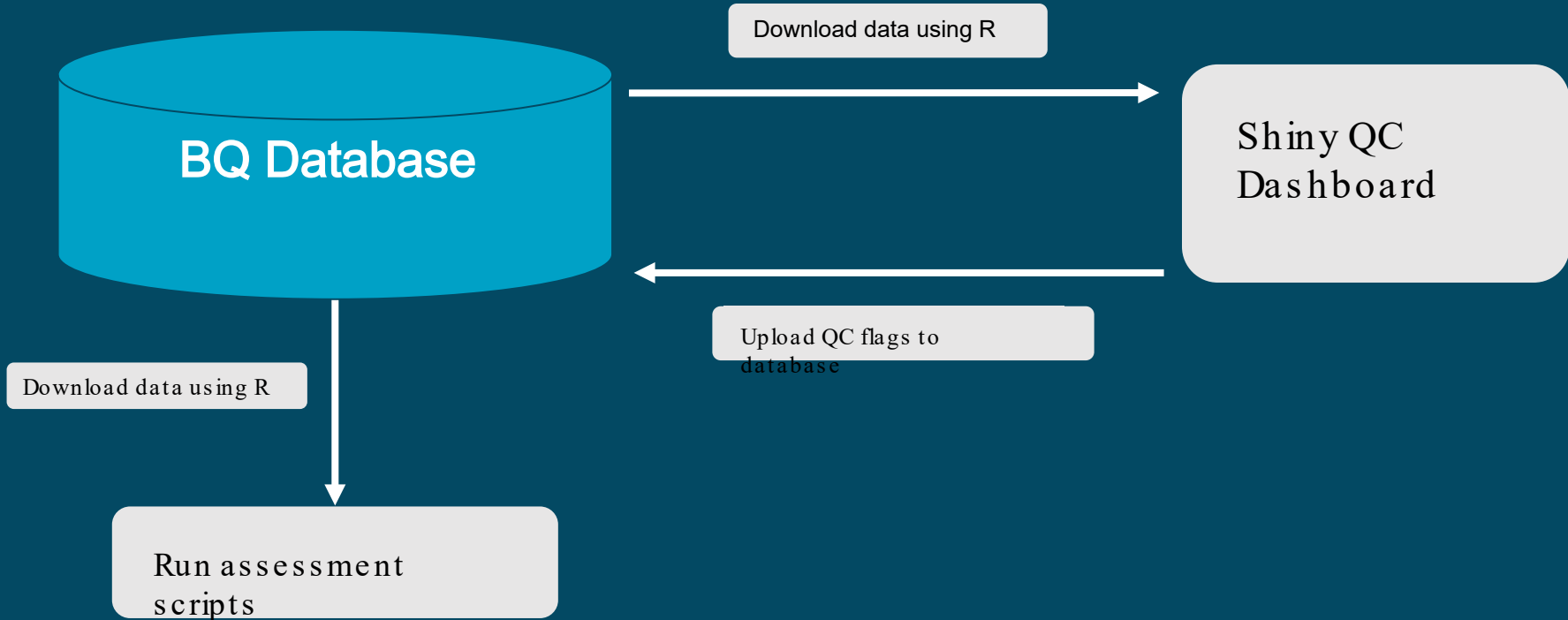
- Access data from:

  - Deployed loggers

  - USGS sites

  - NOAA weather data

- Fast, centralized access supports efficient workflows

# Example Workflow : QC Data

- Workflow:
  - Pull data from DB into RStudio
  - Load into dashboard for visualization
  - Overlay continuous data with contextual data
  - Identify and flag issues (e.g., out - of-water sensors)
  - Push QC flags back into BigQuery
- Visual checks improve confidence in the QA process

# Turning QC Data into Assessment

# Closing the Loop

- BigQuery Database
  a. Fast access
  b. Flexible
  c. Can be complex
- Takeaway: Cloud tools enable access to high-resolution water quality data
- Current Data Volumes:
  a. 1800 datasets (including ~340 USGS sites)
  b. Water Data: 29M rows
  c. Streaming Data: 1M rows
  d. USGS Data: 18M rows
- Looker studio dashboards can be developed with minimal code
- Monthly Cost: ($25-30) - January 1-June 2, 2025 = $127.05