# Exploration of Wisconsin Based Wastewater Epidemiology

Marlin Lee, Kyllan Wunder, Abe Megahed
Data Science Institute, University of Wisconsin-Madison
Special thanks to Dagmara Antkiewicz, Adelaide Roguet
from the Wisconsin State Lab of Hygiene

# Outline:

1. Problem Definition
   - What's WBE(Wastewater-Based Epidemiology)
   - Data sources - where is the data from (collection, analysers, reporters)
   - Data challenges
2. Investigations
   - Normalization
   - Cleaning
   - Analysis
3. Conclusions

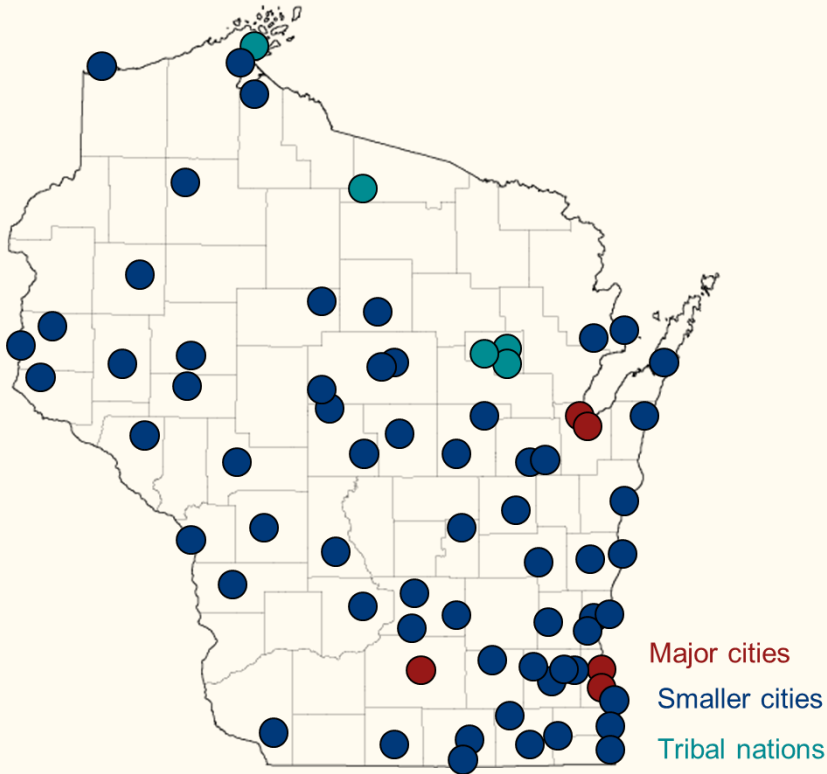# What is Wastewater-Based Epidemiology (WBE)?

Process:

- Collect samples from sewersheds
- Analyze samples at State Lab of Hygiene
- Report results at Department of Health services

Goals:

- Help people make better policy
- Identify outbreaks before they manifest in the case data

# Wisconsin Covid-19 Data Sources



## Wastewater Data

- 82 Sites
- ~3 years
- 1 - 6 days / week
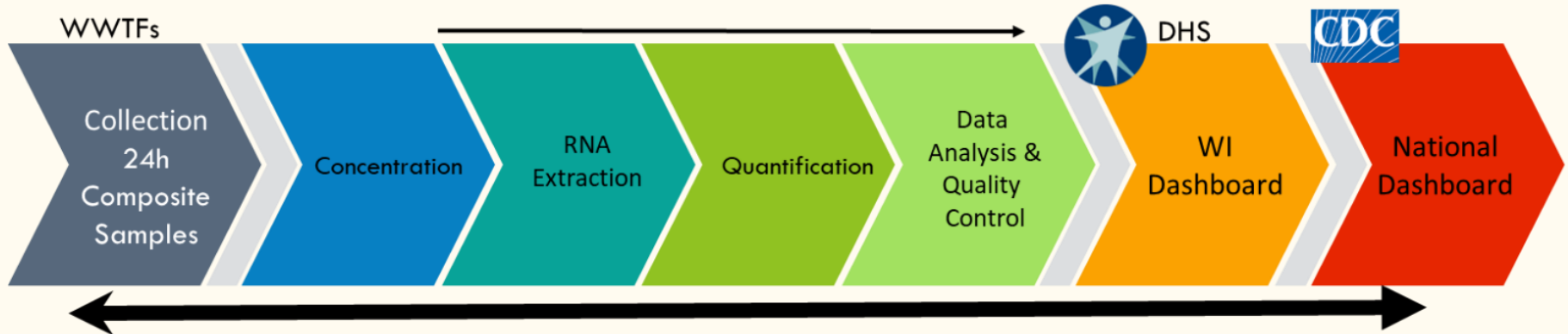- 64,000 samples

## Case Data

- 82 Sites
- ~3 years
- 7 days / week

Major cities

Smaller cities

Tribal nations

# WBE is a Multi-Organization Collaboration
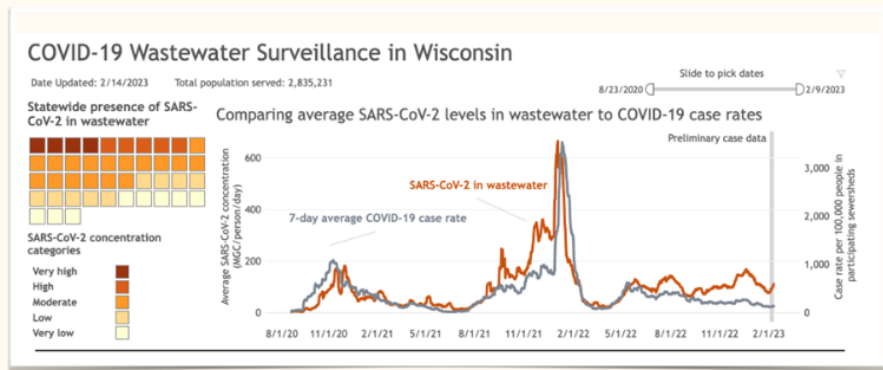
Regional

- DHS
- WSLH
- UM Milwaukee

National

- CDC

# How WBE Information is Communicated to the Public



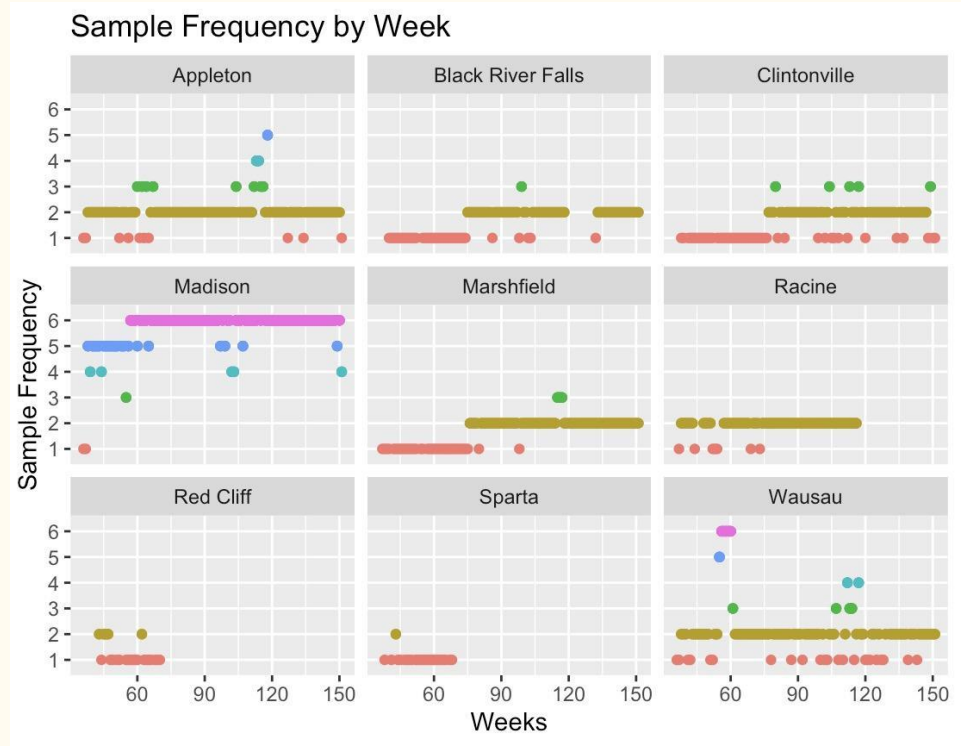WI Department of Health Services (**DHS**) dashboard

https://dhs.wisconsin.gov/covid-19/wastewater.htm



Centers for Disease Control and Prevention (**CDC**) dashboard

https://covid.cdc.gov/covid-data-tracker/#wastewater-surveillance

# Challenges of Working with the Data

- Data is inherently noisy
- Sampling differences (once / week vs. many times / week)
- Analysis method differences (qpcr vs. dpcr)
- Many cofactors (population, time etc.)

# Investigations

- Cleaning
  - Outliers
  - Smoothing

- Normalization

- Analysis
  - Variance analysis
  - Offset analysis



Presented on April 17, 2023 Midwest Chapter SETAC

# Outliers

- Wastewater measurement data tends to be noisy

- Removing outliers allows improved detection of trends in the data



Black River Falls - population: 5000

# Outlier Detection

Solution:

- Look at local spikes
    - Only requires 7 data points / at most 2 week future knowledge
- Pick threshold based on historical spike data
    - Controls for difference frequencies

| site | case freq | waste freq |
|---|---|---|
| Madison | 7 | 6 |
| Marshfield | 7 | 2 |
| Black River Falls | 7 | 2 |

Problems:

- Inconsistent frequency of wastewater collection
    a. Solution must be able to handle this
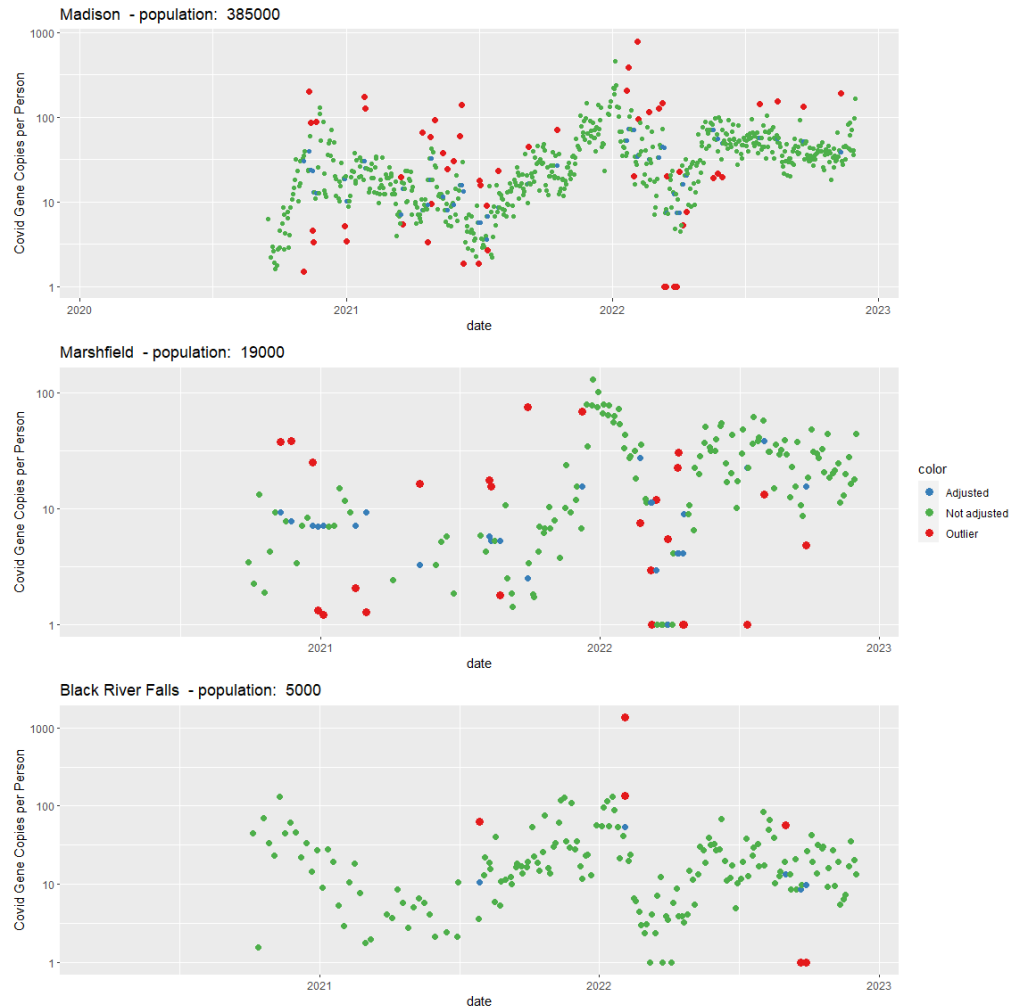- Low minimum frequency
    a. One measurement / week means 2 future data points is a soft cap on influence
    b. Removing data means no measurement for the week

# Outlier Removal

- Removing outliers fixes issues on a small time scale

- Does not meaningfully improve overall trend

# Smoothing

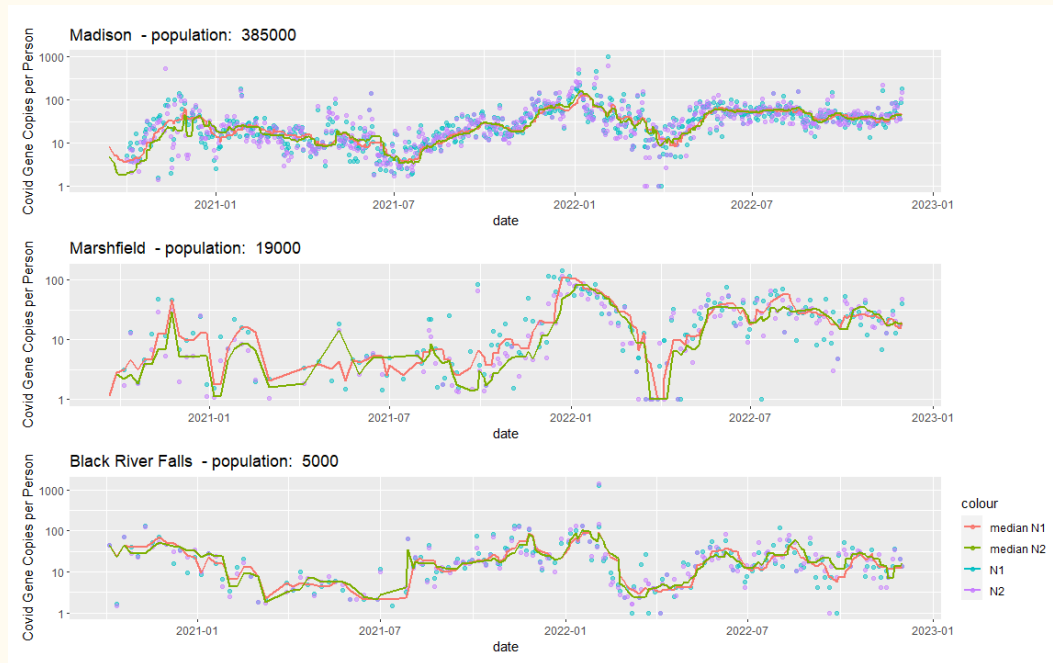- Smoothing removes high frequency noise in the data

Approach:

- Use right aligned Median smoothing

Issues:

- Sitewise Inconsistent frequency of wastewater collection

Results:

- Smoothed data (lines) provide better trending indicators than non-smoothed data (points)

# Normalization

Goals:

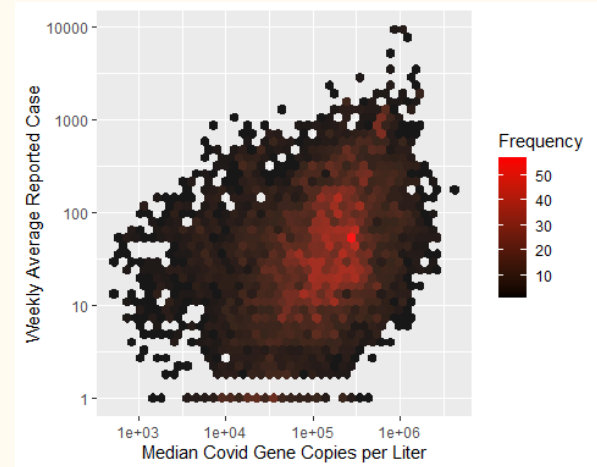- We want measurement to represent the true presence of Covid in the community

Concerns:

- Does the signal scale with population?
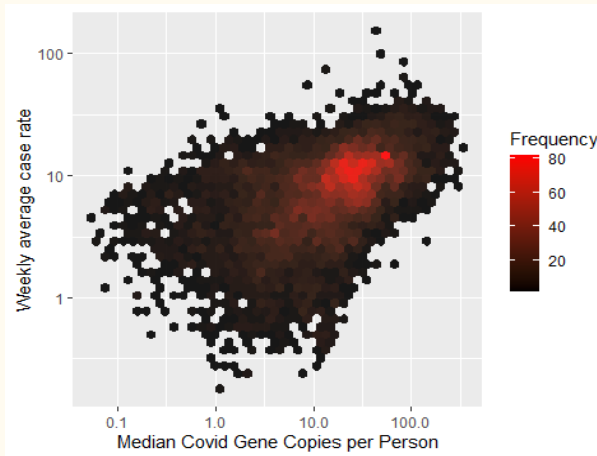- Does the signal scale with the collection method (tests, flow)?

Results:

- Reduced variance
- Better captures underlying signal
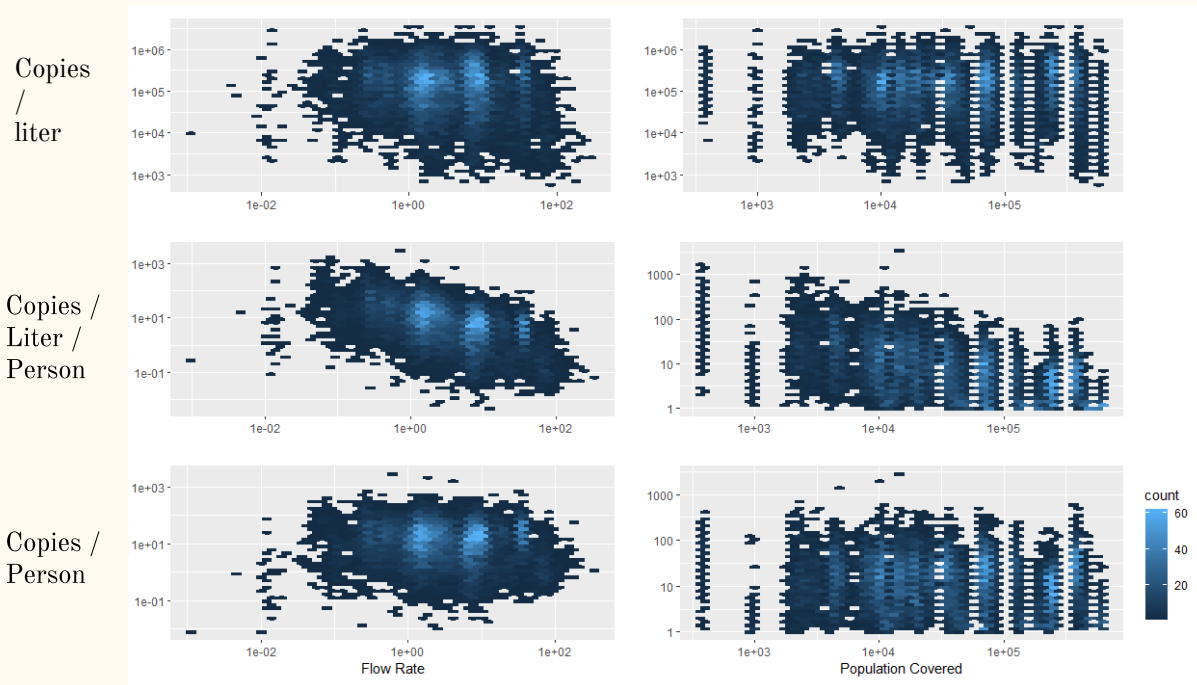- Reduced model error

Before:



After:

# Normalization of Gene Copies

Conclusion:

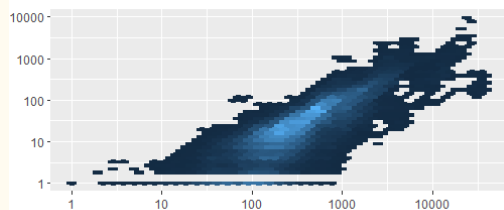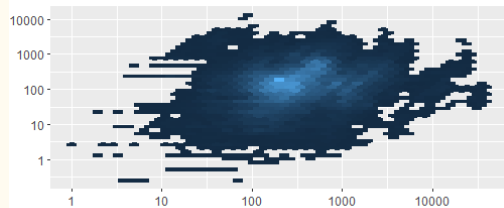Gene * Flow Rate / Population has the lowest covariation

# Normalization of Reported Cases

Conclusion:

Percent positive has the lowest covariation
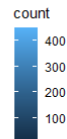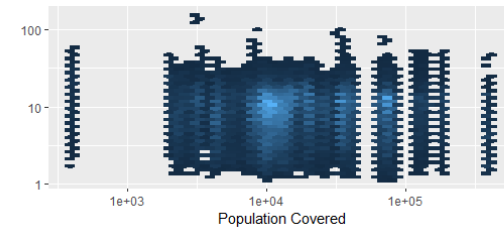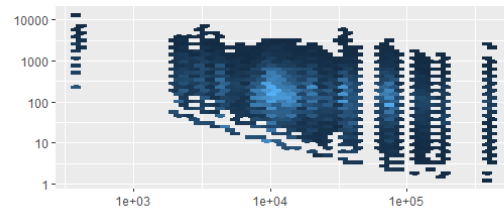
# Basic Level of Detection (LOD) Filtering

- Reported LOD changed over:
  - Time
  - Location
  - Method
- The occurrence of low case rates appears to be correlated with being below the LOD, independent of variations in gene copies.

# Level of Detection Regulating

- Controlling for values below the LOD threshold helps to reveal a more distinct trend.
- There is a decently strong connection between cases and N1/N2 concentration discounting time.

Drop values below LOD



No values removed



Set values below LOD to LOD / 2

# Correlation Between Reported Cases and Gene Copies

- Reported cases (above) track with gene copies (below) over time.

N1, N2 = gene markers

# Variance Analysis Overview

Made easy:

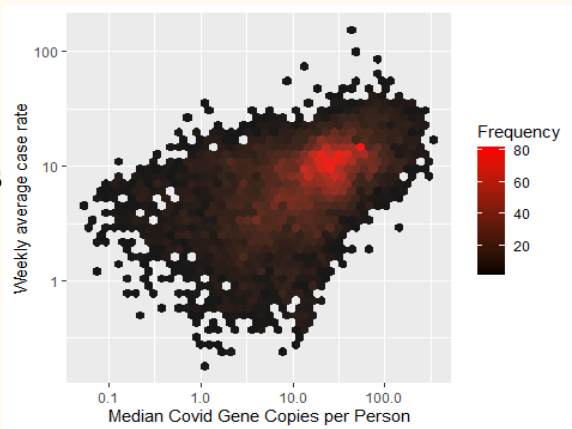- WSLH ran a high frequency testing period at 10 sites for 42 days
- Each day 3 filters were collected with 3 measurements each meaning 9 total measurements
- Used to calculate sources of error in each step of collection

# Variance Analysis Results

- Normalizers (HF183, CrP, and PMMoV) have lower systemic variance than N1 and N2.

- PMMoV has higher systemic variance than other normalizers

- N2 has higher variance at ever level then N1

# Time Series Analysis

- Goal:
  Determine temporal offset between cases and wastewater (N1/N2) measurements.

- Results:
  Future days more important



Correlation of past and future days of N1 and N2 to current day cases

# Offset Analysis Results

- Offset between wastewater and case data varies by variant.

- Offset is sometimes lagging and sometimes leading case data.

- Variant and population size both affect offset.

- Reasons for differences between variants are unknown.



☐ Offset out of range >10 or <-10

# Conclusions

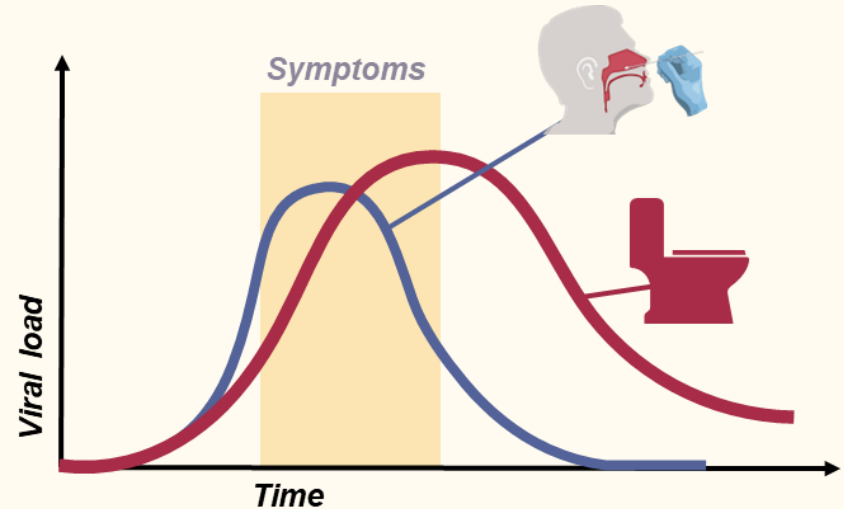- Outlier detection helps fix local (short term) issues but does nothing for the wider trend.
- Smoothing works but creates lagging issues.
- Normalization has huge impact on signal comparison.
- A more comprehensive breakdown of the sources of variance allows for better, more predictive models.

# Future Directions

- Achieve a better understanding of variance in order to create more comprehensive model.
- Improve estimates of shedding and case testing viral load distributions.
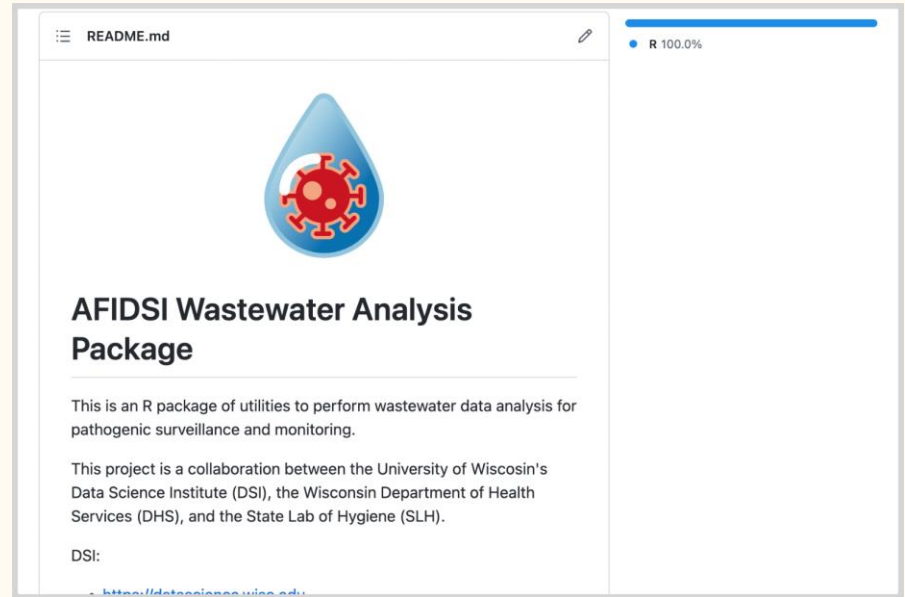- Continue optimizing current methods.

# Final Products

- Presented two poster sessions
- GitHub repository
- Contact us if you'd like more information.

Contact:
Marlin Lee
mrlee6@wisc.edu



https://github.com/AFIDSI/Covid19-Wastewater-Analysis