# Shedding Hub: An Open Science Portal for Existing Pathogen Shedding Data and Models

Yuke (Andrew) Wang, PhD, MSPH [1,2,3]

yuke.wang@emory.edu

Till Hoffmann, PhD, MPhys [4]

thoffmann@hsph.harvard.edu

[1]Hubert Department of Global Health, Rollins School of Public Health, Emory University
[2]Emory Center for Infectious Disease Modeling and Analytics and Training Hub, Emory University
[3]Center of Global Safe Water, Sanitation, and Hygiene, Emory University
[4]Department of Biostatistics, Harvard T.H. Chan School of Public Health, Harvard University

IDM Annual Symposium, October 1st, 2024

# Center for Infectious Disease Modeling and Analytics & Training Hub (CIDMATH)

# InsightNet (National Outbreak Analytics & Disease Modeling Network)



*Illustration of CFA's partners working to detect and control an infectious disease outbreak. https://www.cdc.gov/insight-net/php/about/index.html*

- Established in 2023 by CDC Center for Forecasting and Outbreak Analytics (CFA)
- Focuses on training, analytical tool development, and advancing the analysis and use of data about infectious disease spread
- brings together >100 academic and private partners and health departments

# CIDMATH



*InsightNet's partners. https://www.cdc.gov/insight-net/php/about/index.html*

- 13 centers funded through the CFA
- Emory CIDMATH is a Center of Innovation
- Partners include Georgia DPH, Kaiser Permanente of GA, and the Georgia Emerging Infections Program (EIP)
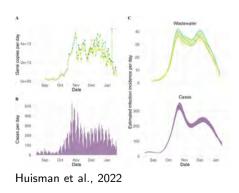
# Motivation

# Wastewater Surveillance



- Wastewater surveillance (WWS) is an approach for monitoring specific pathogen(s) circulating in a population by examining sewage samples.

- Pathogens shed in feces, urine, sputum, and vomit are aggregated in the sewage system.

- A well-designed WWS provides actionable information, including **certification of elimination, early warning, nowcasting/predicting trends, and identification of hotspots**.

# Sources of Uncertainty for WWS



Wade et al., 2022

# Wastewater-based Epidemiology



Huisman et al., 2022

- The interpretation of WWS results in terms of clinical cases depends on quantitatively well-characterized shedding information for pathogens and biomarkers



Phan et al., 2023

# Current Knowledge of Shedding

- Most shedding data have been collected in clinical studies or human challenging studies
- Limited raw shedding data are openly available
- Shedding data are not standardized
- No public accessible tutorial or feasible tool for modeling shedding dynamics
- No community portal for learning and contributing shedding knowledge

Shedding Hub

# Open Science



Bertram et al., 2023

Principles of Open Science:

- Accessible
- Verfiable
- Reliable
- Reproducible
- Sustainable

# Shedding Hub Organization

# Shedding Hub Data Governance Framework



This Data Governance Framework is an extension of the tillahoffmann/shedding repository.

## Data Structure



- Includes raw data files, markdown files to process data, standardized data in YAML format, and a schema file
- Public **accessible** on GitHub

## Data Processing



- **Verifiable** with information sources
- **Reproducible** with open code

# Data Standardization



- JSON schema (json-schema.org) is a standard to specify the structure of data, validate it, and include documentation about each field
- All data uploaded will be validated against the schema automatically

# Data Validation



To generate **reliable** data:

- Data will be checked and reviewed by at least two reviewers
- Functions were created to generate automatic data summary
- Conversations and decisions will be documented on GitHub

# Modeling Tutorial and Code



- Tutorial for modeling shedding data in Python/R
- **Reproducible** using markdown files

## Contribution



- **Sustainable** with partnerships and contributions from the research community

# Shedding Hub Website



shedding-hub.github.io

# Work Plan

| Phase 1 | Phase 2 | Phase 3 |
|---|---|---|
| **Objectives**: Build a Shedding Hub Team; Develop the Basic Structure and Workflow | **Objectives**: Expand to Additional Pathogens and Biomarkers; Invite Contribution of Data and Models | **Objectives**: Develop analytical tools (e.g., dashboard, packages) and promote usage and contribution |
| **Data Sources**: Open Access Data, Published in Literature **Priority**: Quantifiable Longitudinal Measurements | **Data Sources**: Limited Access Data **Priority**: Quantifiable Longitudinal Measurements | **Priority**: Semi-quantifiable (CT values) or Non-quantifiable (Presence/Absence) Measurements |
| **Models**: Developed within the Shedding Hub Team | **Models**: Published Models by Experts | **Models**: Contributed Models by Research Community |
| **Prioritized Biomakers**: Pathogens of Interest for WWS | **Prioritized Biomakers**: Additional Pathogens and Biomarkers | **Prioritized Biomakers**: Additional Pathogens and Biomarkers |
| *May 2024 – Sep 2025* | *Oct 2025 – Sep 2026* | *Oct 2026 – Sep 2027* |

## Use Cases

- To support wastewater-based epidemiology applications estimating disease incidence for various infectious diseases
- To better understand of sensitivity of different disease diagnostic methods (nasopharyngeal swab vs. rectal swab)
- To support decision making for disease control and prevention policys (e.g., how long the quanrantine period should be?)
- To support wastewater monitoring for drug use

# Acknowledgements

## Shedding Hub Team



Yuke Wang

Till Hoffmann

Weifei Xiao

Youwei Hu

Zirui Chen

Haisu Zhang

## CIDMATH Team