# Data Modernization: Improving the usefulness of genomic data

Kelsey Florek, PhD, MPH
Senior Genomics and Data Scientist
Wisconsin State Laboratory of Hygiene
May 21, 2024

Slides live at:

www.k-florek.net/talks

Wisconsin State
Laboratory of Hygiene
UNIVERSITY OF WISCONSIN–MADISON

# Supported By



AWS Diagnostic Development Initiative (DDI)

1. **Necessities of Next Generation Sequencing Capacity Building**

2. Blueprints for an NGS Data Solution

3. Simplifying Genomics for Public Health Partners

# Expanding Genomic Sequencing Capacity

<span style="color:red">**Pre SARS-CoV-2 Pandemic**</span>

- 4x Illumina MiSeq

- 1x ONT MinION

# Expanding Genomic Sequencing Capacity

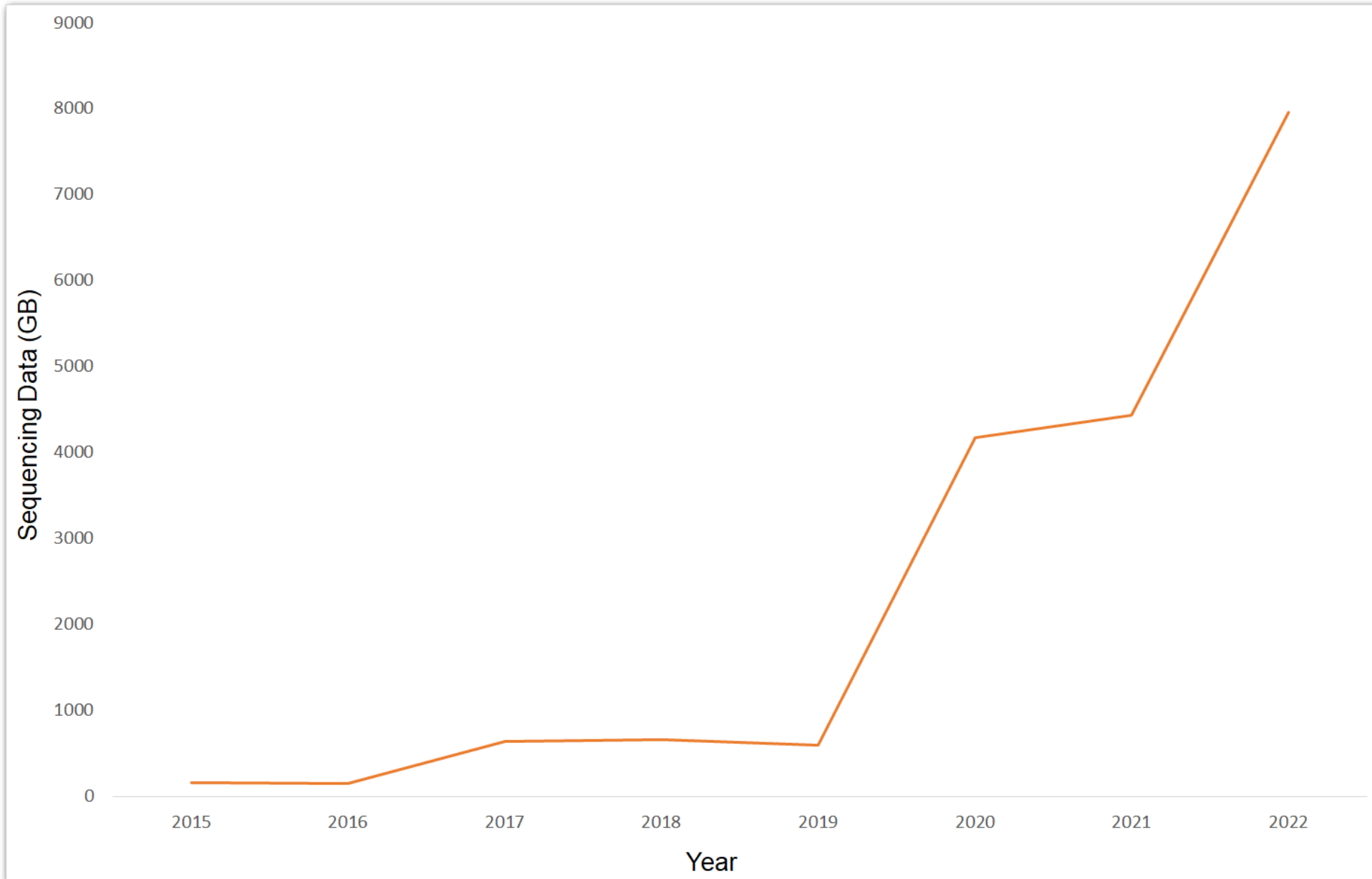## Pre SARS-CoV-2 Pandemic

- 4x Illumina MiSeq

- 1x ONT MinION

## Post SARS-CoV-2 Pandemic

- 4x Illumina MiSeq

- 2x NextSeq 2000

- 1x ONT GridION

- 1x Eppendorf epMotion

- 1x Tecan Fluent 780 NGS Dream Prep

# Expanding Genomic Sequencing Capacity

## Pre SARS-CoV-2 Pandemic

- 4x Illumina MiSeq

- 1x ONT MinION

## Post SARS-CoV-2 Pandemic

- 4x Illumina MiSeq

- 2x NextSeq 2000

- 1x ONT GridION

- 1x Eppendorf epMotion

- 1x Tecan Fluent 780 NGS Dream Prep

# Over 900% increase in sequencing data generation capacity

NGS Data Storage

# Improvements in Analytical Approaches

## Old Approach

- Entirely Python Based

- Limited logging and fault tolerance

- Required installing complicated and often conflicting dependencies

# Improvements in Analytical Approaches

## Old Approach

- Entirely Python Based

- Limited logging and fault tolerance

- Required installing complicated and often conflicting dependencies

## New Approach

- Nextflow - Nf-Core Based

- Containerized Steps

- Detailed Logging

- Compatible with a variety of Cloud and HPC environments

- Supports a high degree of job parallelization and horizontal scalability

7

# Bioinformatics analytical infrastructure

- Highly scalable and capable of managing burst data

- Highly reliable and fault tolerant

- Cost effective

- Adaptable to changing needs

- Detailed logging and traceability

# WSLH Bioinformatics Analytical Infrastructure

**AWS Batch**

AWS Batch automatically provisions compute resources and optimizes the workload distribution based on the quantity and scale of the workloads.

# nextflow *tower*

Nextflow Tower is an intuitive centralized command post that enables data analysis at scale. With Tower, users can easily launch, manage, and monitor scalable Nextflow data analysis pipelines and compute environments on-premises or across the cloud providers of their choice.

12

# Nextflow Tower - Dashboard

# Nextflow Tower - Monitor

Seqera Labs - Nextflow Tower

# Connecting Data Across Siloed Systems

1. ~~Necessities of Next Generation Sequencing Capacity Building~~

2. ~~Blueprints for an NGS Data Solution~~

3. **Simplifying Genomics for Public Health Partners**

# Need for a centralized resource

## AMD Bioinformatics Regional Resource - Midwest Region



Legend:
- Southeast
- Northeast
- Central
- Mid-Atlantic
- Midwest
- Mountain
- West

Territories: AS GU PR VI MP FM PW MH

**Ad-hoc Analytical Support**

**Provision of Computational Resources**

# Need for a centralized resource

## SARS-CoV-2 Genomic Surveillance

WSLH SARS-CoV-2 Dashboard

# COVID-19 Genomics UK (COG-UK) CLIMB-COVID



Nicholls et al. 2021 *Genome Biology*
StaPH-B Monthly Webinar Oct 2021

20

# COVID-19 Genomics UK (COG-UK) CLIMB-COVID

# Easy Genomics Partnership

**Two Bulls/DEPT®**

Digital health product development with care.

**Amazon Web Services**

On-demand cloud computing web services.

**Wisconsin State Laboratory of Hygiene**

Wisconsin's Public, Environmental and Occupational Health Laboratory Since 1903

# Easy Genomics - Minimal Viable Product

- Simplify the process of launching and monitoring workflows

- Provide the ability for users to upload sequence data through the web browser

- Allow users to download analysis results through the web browser

- User/Lab separation

# Easy Genomics - Sequence Data Upload

# Easy Genomics - Launch

Workflows    Sequencing    Pipelines    Users    KF

← Back

## Launch Pipeline
Launch a new pipeline

**Step 1**
Select Blueprint

**Step 2**
Required Params

**Step 3**
Confirm launch details

**Workflow Name**

irreverent_cori

A unique name for the workflow.

**Select Pipeline**

WSLH-EG-viralrecon                                              ⌄

Select a pipeline to run

Next

# Easy Genomics - Launch



Easy Genomics

Workflows    Sequencing    Pipelines    Users    KF

← Back

## Launch Pipeline
Launch a new pipeline

**Step 1**
Select Blueprint

**Step 2**
Required Params

**Step 3**
Confirm launch details

### Input/output options
Define where the pipeline should find input data and save output data.

**input**

| Choose File | test_sample_sheet.csv |

s3://dev-wslh-sequencing-inbox/uploads/test_sample_sheet.csv

Path to comma-separated file containing information about the samples you would like to analyse.

**platform**

illumina

NGS platform used to sequence the samples.

**protocol**

amplicon

Specifies the type of protocol used for sequencing.

**outdir**

s3://dev-wslh-sequencing-analyses

The output directory where the results will be saved. You have to use absolute paths to storage on Cloud infrastructure.

**email**

email

Email address for completion summary.

# Easy Genomics - Monitor

# Easy Genomics - Monitor

Easy Genomics

Workflows    Sequencing    Pipelines    Users    KF

← Back

## wise_pauling
View a detailed look at this workflow

Job Details    Job Stats

### Tasks Stats

| | Pending | | Submitted | | Running |
|---|---|---|---|---|---|
| | 0 Tasks | | 24 Tasks | | 1 Tasks |

| | Cached | | Succeeded | | Failed |
|---|---|---|---|---|---|
| | 0 Tasks | | 52 Tasks | | 0 Tasks |

### Aggregate Stats

| | Wall Time | | CPU Time | | Total Memory |
|---|---|---|---|---|---|
| | 0 seconds | | 12.0 CPU hours | | 141.41 GB |

| | Disk Read | | Disk Write | | Estimated Cost |
|---|---|---|---|---|---|
| | 194.96 GB | | 85.8 GB | | $0.460 |

Cancel Workflow

# Easy Genomics - Roadmap

# Easy Genomics - Roadmap

- 2024 Spring - Deploy Easy Genomics for internal use

# Easy Genomics - Roadmap

- 2024 Spring - Deploy Easy Genomics for internal use

- 2024 Early Summer - Open Access to SARS-CoV-2 Sequencing Laboratories

# Easy Genomics - Roadmap

- 2024 Spring - Deploy Easy Genomics for internal use

- 2024 Early Summer - Open Access to SARS-CoV-2 Sequencing Laboratories

- 2024 Mid Summer - Easy Genomics MVP Update

# Acknowledgments

Abigail Shockey, PhD

Christopher Jossart, MPH

Dustin Lyfoung, MS

Thomas Blader

Eva Gunawan, MS

## Special Thanks

- UW-Madison Public Cloud Team

- UW-Madison Office of Cybersecurity