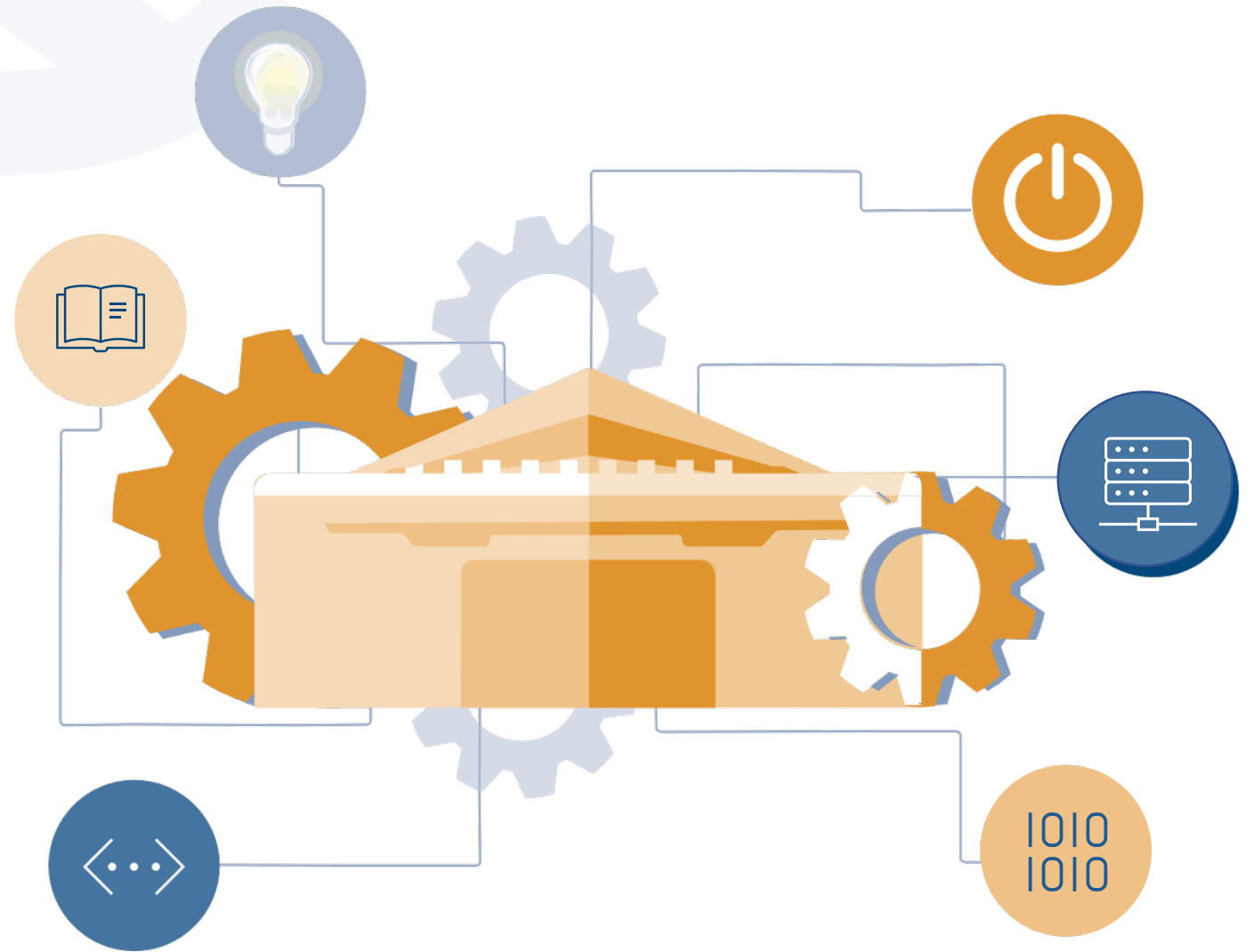


Current State of Infrastructure and Interoperability of Pathogen Genomics Data in the U.S.

Steve Sherry, Ph.D.
Acting Director, National Library of Medicine, NIH



Public sequence data supports pathogen surveillance and outbreak detection

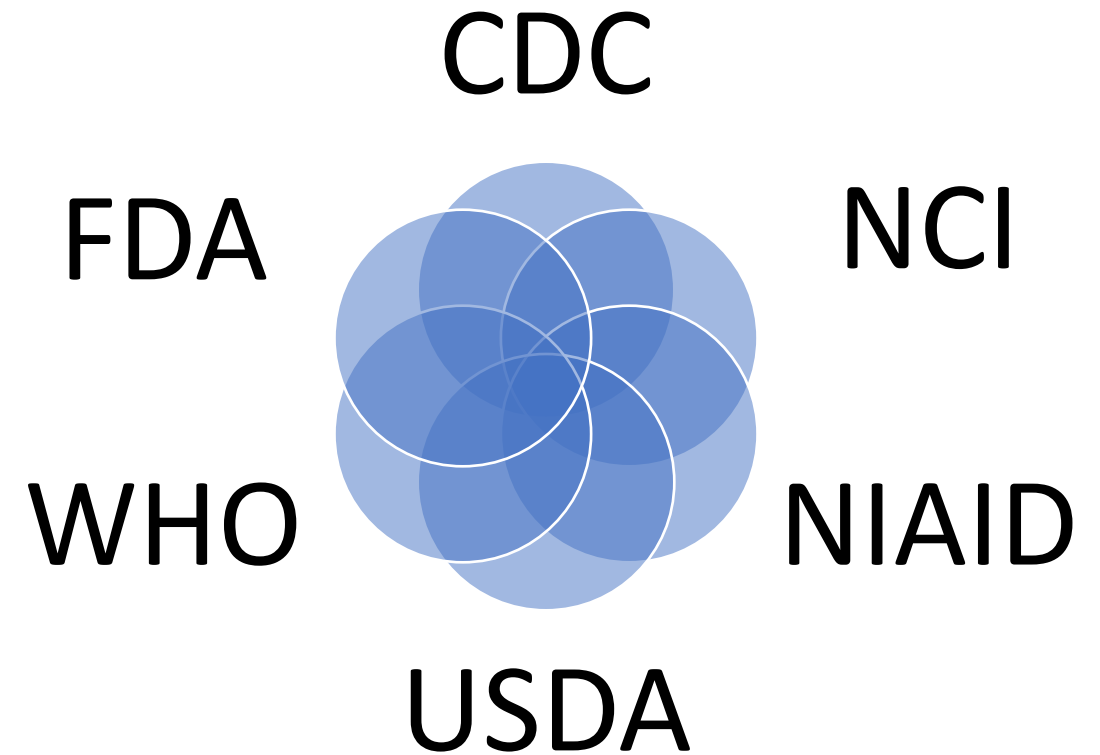
- Rapidly identify sources of outbreak to target interventions and remove threats
 - Foodborne illness
 - Contaminated medical devices
- Identify common sources of infection and target initiatives aimed at reduction
- Detecting emerging pathogen threats
 - Identify novel types or species that are beginning to cause disease
 - Identify antimicrobial resistance

Public sequence data supports pathogen surveillance and outbreak detection

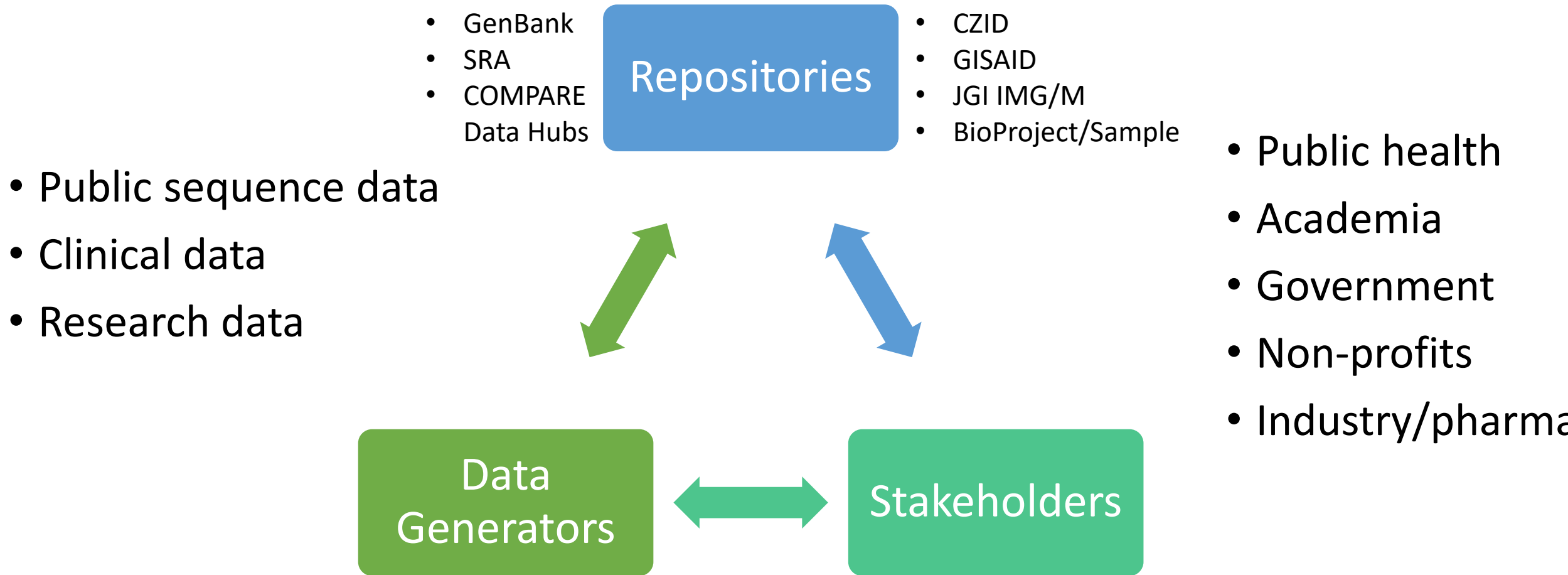
- Can accelerate research across multiple disciplines such as developing interventions or vaccines
- Monitoring ongoing pathogen evolution
 - Evolving resistance to vaccines
 - Evolving genetic resistance to therapeutics such as antibacterial or antiviral agents
 - Understanding pathogen evolution
 - what are the genetic components of virulence and pathogenicity in the pathogen that cause disease
 - What are the host or environmental factors that contribute to the spread of disease

Pathogen Genomics Priorities: Surveillance vs Outbreak

- Multiple groups maintain different lists of high-priority pathogens
 - Dynamic
 - Reflect different interests
- Prioritized pathogens can have different infrastructure and data management requirements:
 - chronic, ongoing relationship associated with monitoring
 - Shorter term, outbreak-related

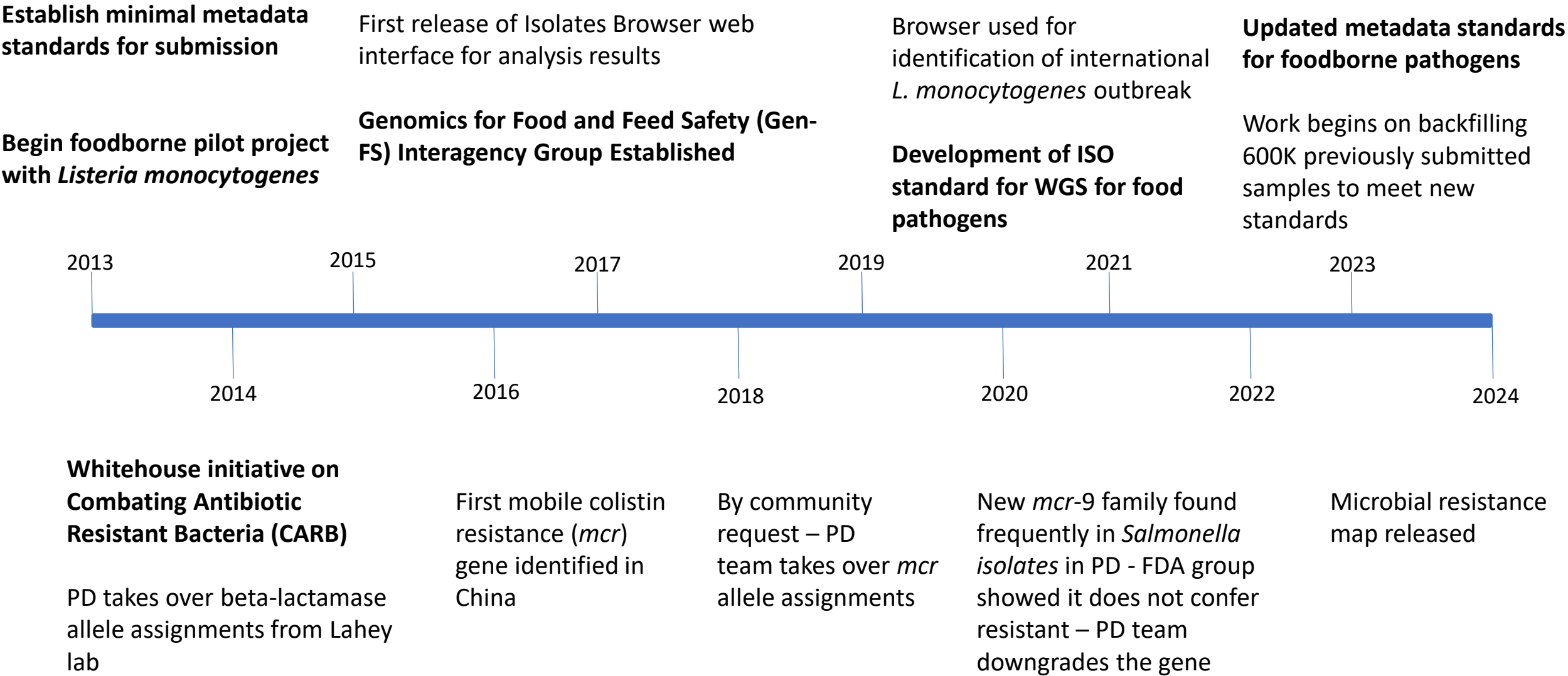


Developing a Pathogen Genomics Data Infrastructure



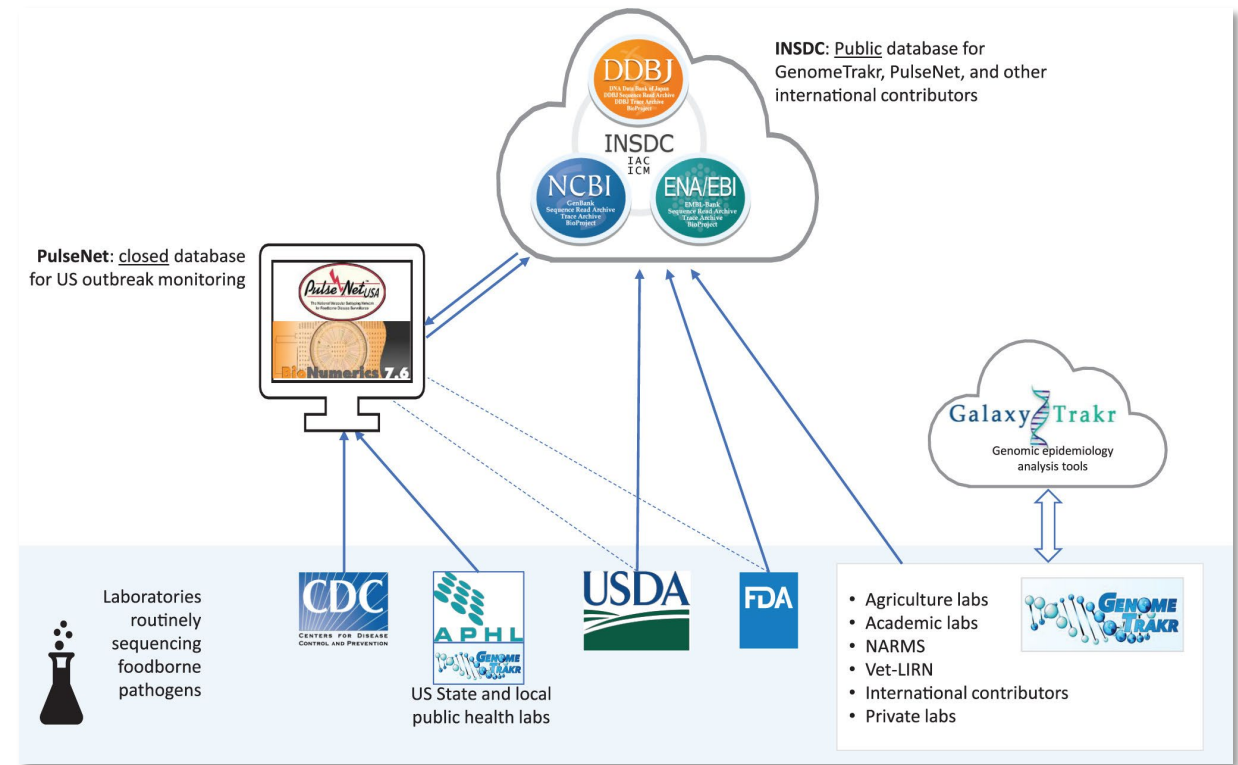
Pathogen Detection (PD) Development Timeline

Funding support via NIAID from CARB



Interagency Collaboration for Genomics for Food and Feed Safety (Gen-FS)

- National Institutes of Health (NIH)
- Centers for Disease Control and Prevention (CDC)
- U.S. Food and Drug Administration (FDA)
- U.S. Department of Agriculture (USDA)
 - Food Safety and Inspection Service (FSIS)
 - Agricultural Research Service (ARS)
 - Animal and Plant Health Inspection Service (APHIS)
- Standardization of isolate metadata for submission
- Harmonized sequencing protocols across CDC PulseNet & FDA GenomeTrakr
- Harmonized quality assurance measures & accompanying quality control checks
- Ensures all WGS data generated by Gen-FS members in the Pathogen Detection database at NCBI meet the minimum quality standards

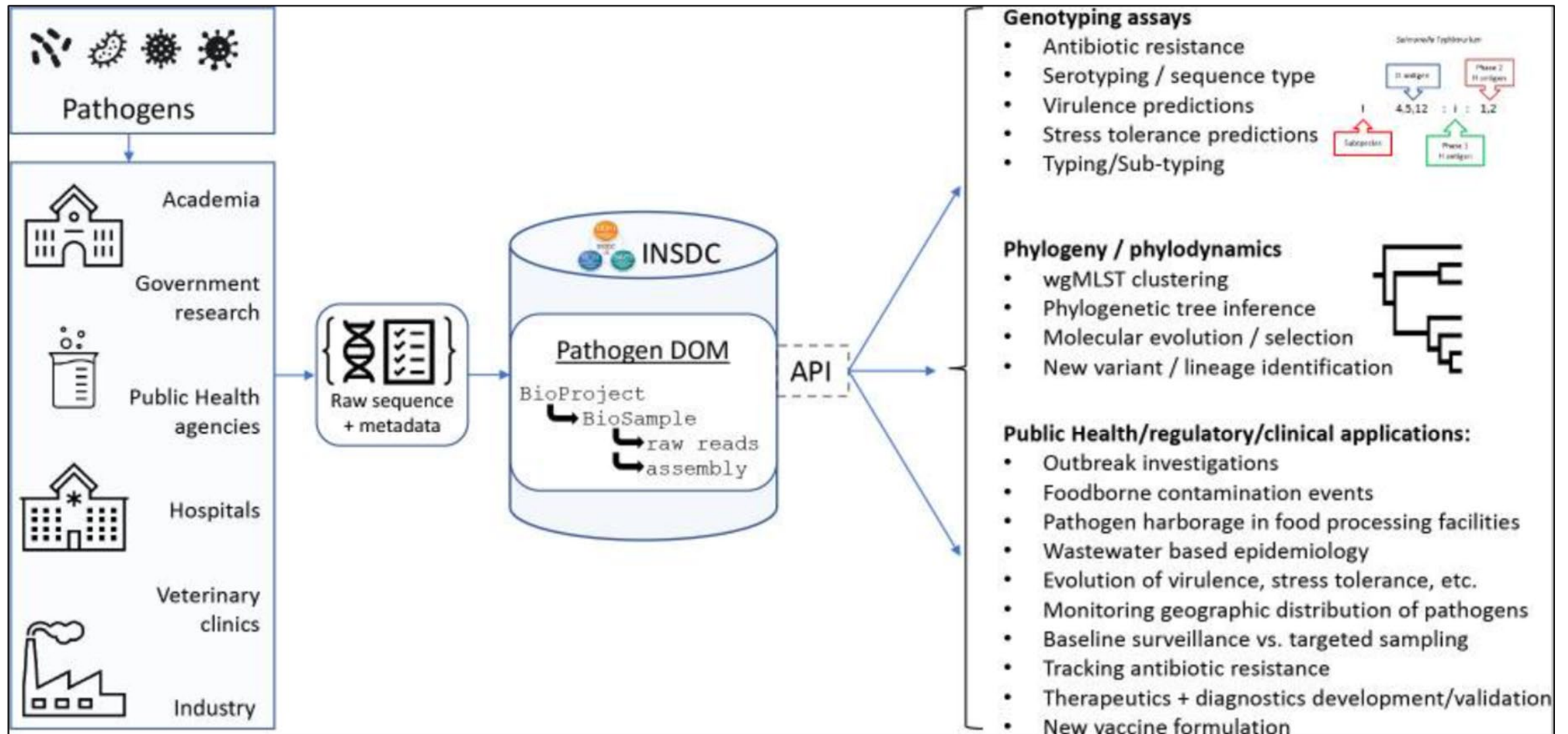


PHA4GE

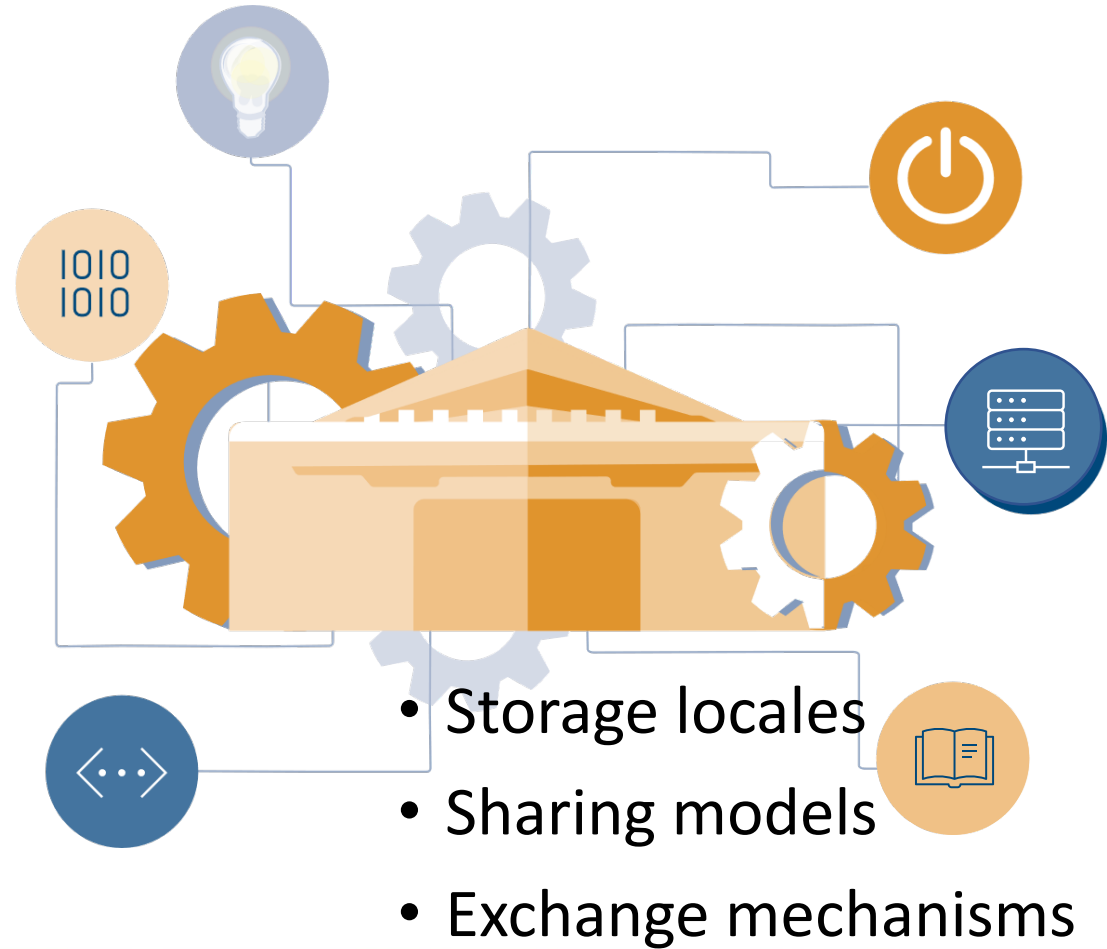
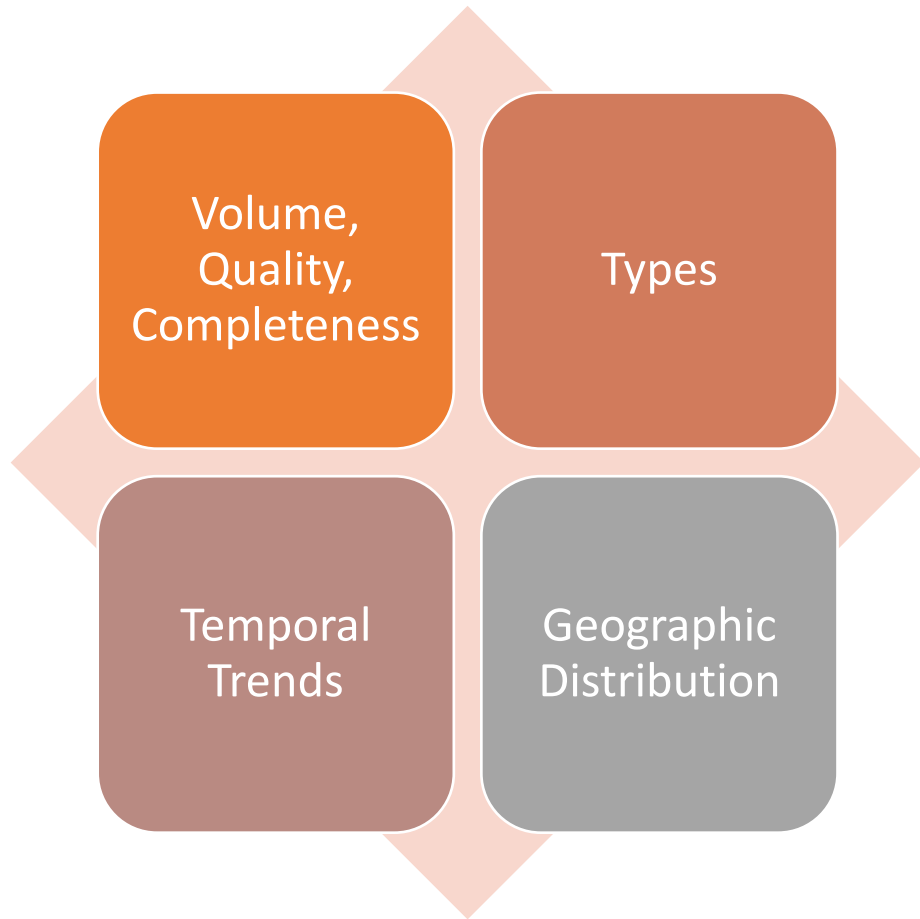
- PHA4GE- a collaborative community working to:
 - Advance the use of open data
 - Empower laboratories to analyze and govern their own data
 - Encourage standardized data structures and interchange formats



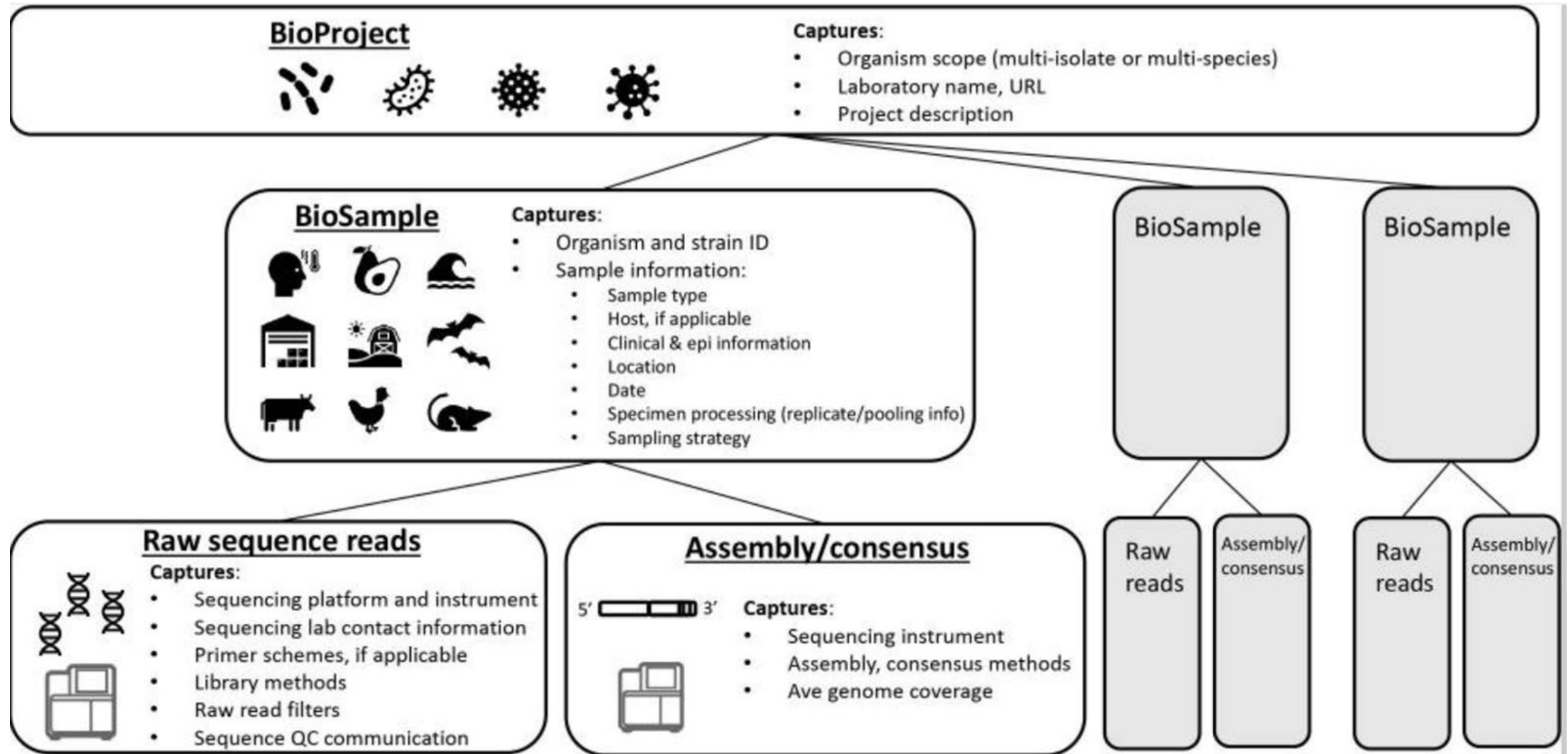
The INSDC compliant Pathogen Data Object Model to structure genomic data for public health applications



Variables in PG Data Management and Reuse



The INSDC compliant Pathogen Data Object Model to structure genomic data for public health applications



One Health Enteric Package

BioSample package developed under Gen-FS metadata working group captures the full One Health sample space for enteric microbes

One Health Enteric package scope



Developed by international group of experts



CORE attributes

- Isolate identifiers
- Collected by
- Date of collection
- Geographic location
- Sampling purpose
- Sampling device
- Project name
- IFSAC category
- Source type
- sequenced by



Human/animal host

- Host
- Host disease
- Host sex + age
- Host tissue sampled
- Animal environment
- Antimicrobials in food
- Animal housing system



Food samples

- Geographic origin
- Intended consumer
- Collection site description
- Food product type
- Label claims
- Food source
- Food processing types
- Food preservation process
- Food additives
- Food contact surface
- Food container wrapping
- Food container integrity



Food facility

- Facility type
- Building setting
- Food processed
- Facility location
- Monitoring zone
- Indoor sampling surface
- Surface material
- Surface material cond.
- Surface orientation
- Surface temperature
- Biocide used
- Animal intrusion



Farm and Environment

- ENVO triad
- Farm type
- Plant growth medium
- watering method
- Relative loc of sample
- Fertilizer administration
- Food cleaning process
- Sanitizer used
- Farm equip. used
- Water samples
- Extreme weather event
- Mechanical damage

Consists of core attributes and optional attributes depending on source

Aim is to better provide outbreak and traceback investigations of bacterial pathogens

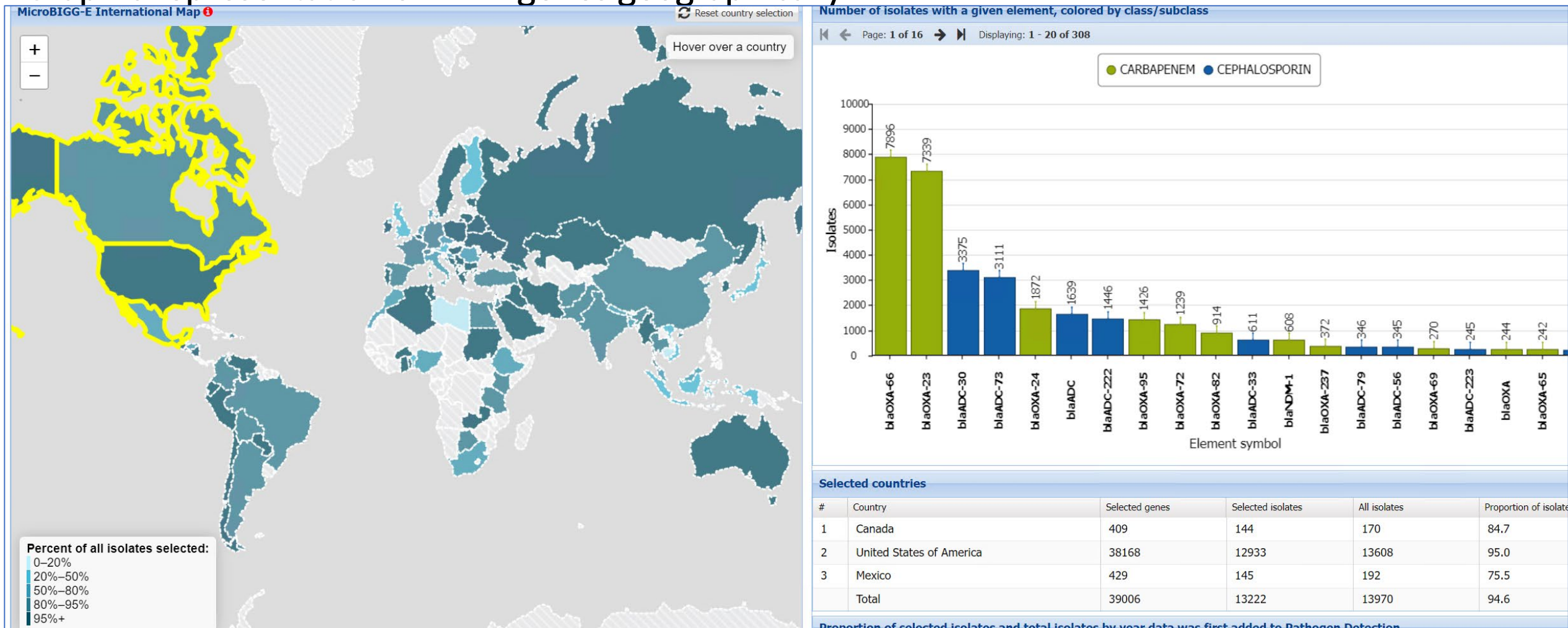
It includes using ontologies for certain attributes to provide better structure for comparison than free text

12

Selected submissions to Pathogen Detection for past year

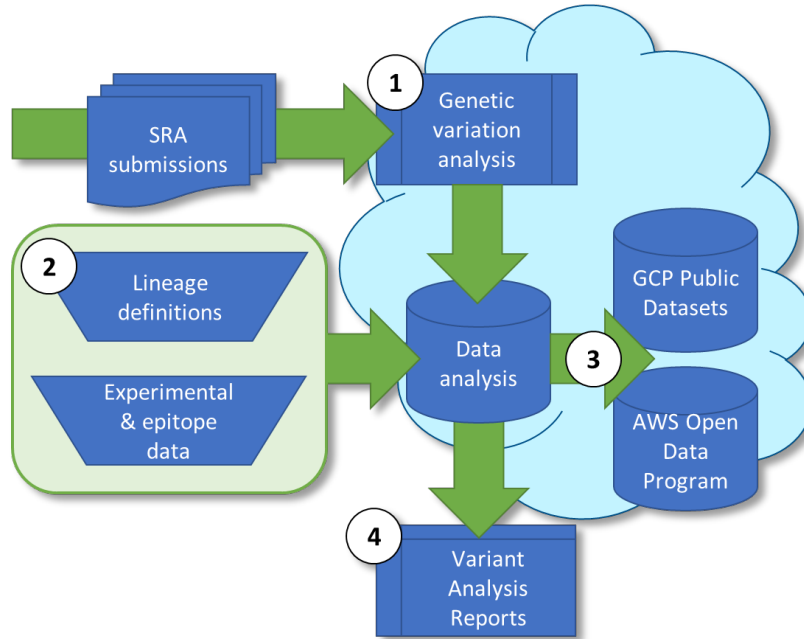
Category	Organism	No. of Isolates	No. of Countries	% Missing Location
Foodborne	<i>Salmonella enterica</i>	95, 433	32	18.95
	<i>Escherichia coli</i> + <i>Shigella</i> spp.	73, 319	94	15.23
	<i>Campylobacter</i> spp. (<i>jejuni</i> and <i>coli</i>)	15, 999	30	11.61
	<i>Listeria monocytogenes</i>	8, 293	12	9.73
	<i>Cronobacter sakazakii</i>	945	12	14.71
Clinical	<i>Acinetobacter baumannii</i>	14, 555	56	0.99
	<i>Klebsiella pneumoniae</i>	19, 224	77	5.09
	<i>Staphylococcus aureus</i>	16, 740	62	2.80
	<i>Streptococcus pneumoniae</i>	31, 037	41	51.93

Graphic representation of AMR genes geographically



Carbapenem and cephalosporin resistance in *Acinetobacter baumannii* in North America from data submitted by end of 2023

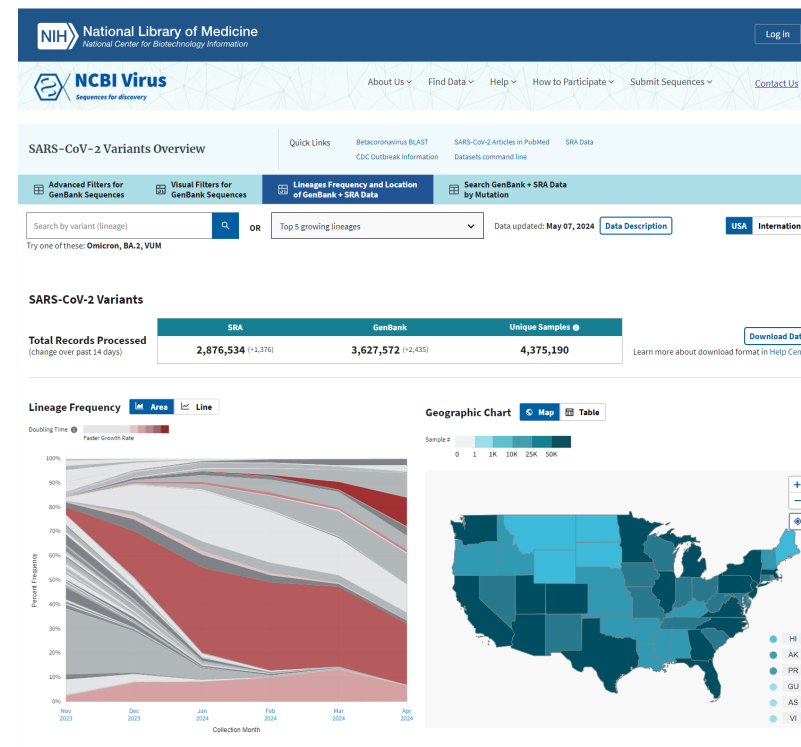
Support of viral pathogen genomics



SARS-CoV-2 variation analysis pipeline

Connor, R, et al. *Viruses*. 2024 Mar 11;16(3):430.
doi: 10.3390/v16030430.

SARS-CoV-2 Variants Overview



The screenshot shows the 'Influenza Virus Data Hub' interface. It includes a search bar, filters for virus type (e.g., AlphaInfluenzavirus, BetaInfluenzavirus), and a table of virus sequences. The table shows the following data:

Accession	Organism Name	Submitters	Organization	Release Date	Isolate
NC_026431	Influenza A virus (A/Calif...	Garten,R.J., et al.	National Center for Biotec...	2015-02-23	
NC_026432	Influenza A virus (A/Calif...	Garten,R.J., et al.	National Center for Biotec...	2015-02-23	
NC_026433	Influenza A virus (A/Calif...	Garten,R.J., et al.	National Center for Biotec...	2015-02-23	
NC_026434	Influenza A virus (A/Calif...	Garten,R.J., et al.	National Center for Biotec...	2015-02-23	
NC_026435	Influenza A virus (A/Calif...	Shu,H., et al.	National Center for Biotec...	2015-02-23	
NC_026436	Influenza A virus (A/Calif...	Garten,R.J., et al.	National Center for Biotec...	2015-02-23	
NC_026437	Influenza A virus (A/Calif...	Garten,R.J., et al.	National Center for Biotec...	2015-02-23	
NC_026438	Influenza A virus (A/Calif...	Garten,R.J., et al.	National Center for Biotec...	2015-02-23	
NC_026439	Influenza A virus (A/Calif...	Garten,R.J., et al.	National Center for Biotec...	2015-02-23	
NC_026440	Influenza A virus (A/Calif...	Garten,R.J., et al.	National Center for Biotec...	2015-02-23	
NC_026441	Influenza A virus (A/Calif...	Garten,R.J., et al.	National Center for Biotec...	2015-02-23	
NC_026442	Influenza A virus (A/Calif...	Garten,R.J., et al.	National Center for Biotec...	2015-02-23	
NC_026443	Influenza A virus (A/Calif...	Garten,R.J., et al.	National Center for Biotec...	2015-02-23	
NC_026444	Influenza A virus (A/Calif...	Garten,R.J., et al.	National Center for Biotec...	2015-02-23	
NC_026445	Influenza A virus (A/Calif...	Garten,R.J., et al.	National Center for Biotec...	2015-02-23	
NC_026446	Influenza A virus (A/Calif...	Garten,R.J., et al.	National Center for Biotec...	2015-02-23	
NC_026447	Influenza A virus (A/Calif...	Garten,R.J., et al.	National Center for Biotec...	2015-02-23	

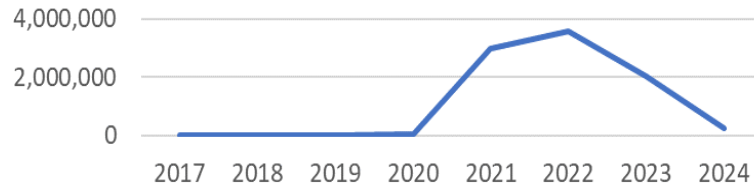
Influenza Virus Data Hub

Selected viral genomic submissions to INSDC in the past year

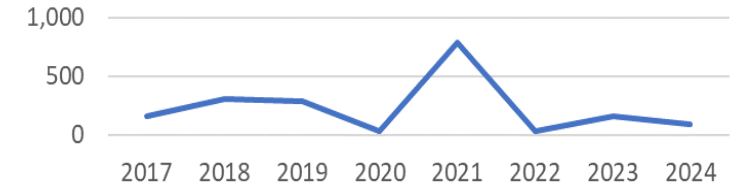
Virus	GenBank Records	No. of Countries	% Missing Location	SRA Records	No. of Countries	% Missing Location
SARS-CoV-2	1,080,335	58	0.0%	428,591	59	0.3%
Influenza A virus	144,282	57	2.0%	6,370	20	4.4%
Influenza B virus	27,448	16	0.6%	368	6	0.0%
Dengue virus	5,416	72	0.0%	1,610	40	0.2%
African swine fever virus	709	35	0.4%	63	10	20.6%
Measles virus	1,104	25	0.0%	54	3	5.6%
Ebolaviruses	207	2	0.0%	14	2	0.0%
Marburg virus	1	1	0.0%	8	2	0.0%
Rift Valley fever virus	101	8	1.0%	68	3	0.0%
Zika virus	141	11	6.0%	674	4	15.4%
HIV-1	51,391	49	12.2%	1,590	11	9.1%

Variability in Submissions of Selected Viral Organisms

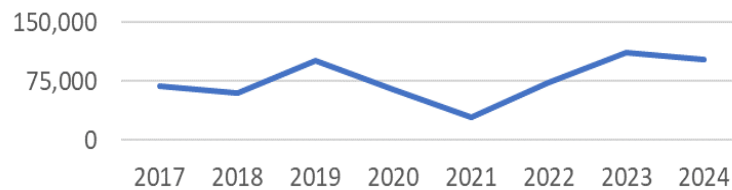
SARS-CoV-2



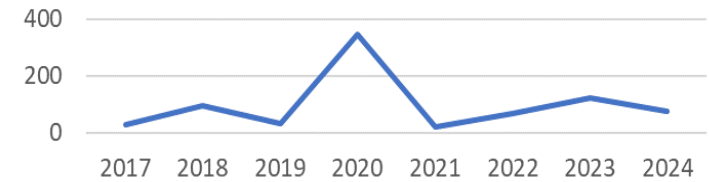
Ebolaviruses



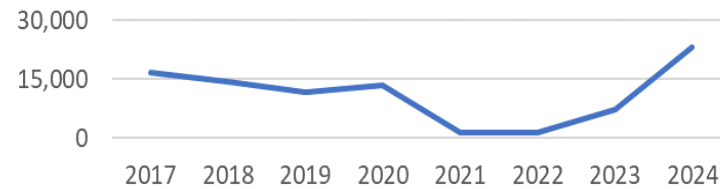
Influenza A virus



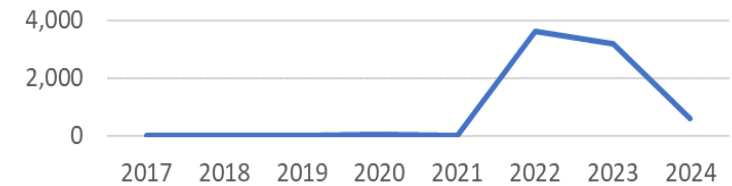
Rift Valley fever virus



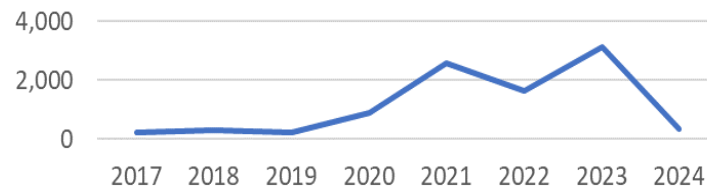
Influenza B virus



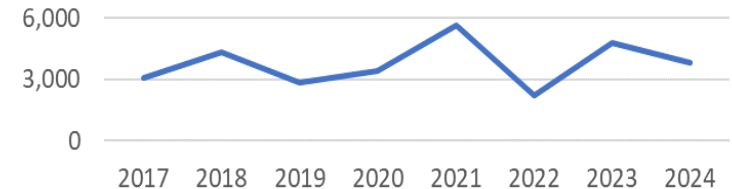
MPOX virus



African swine fever virus



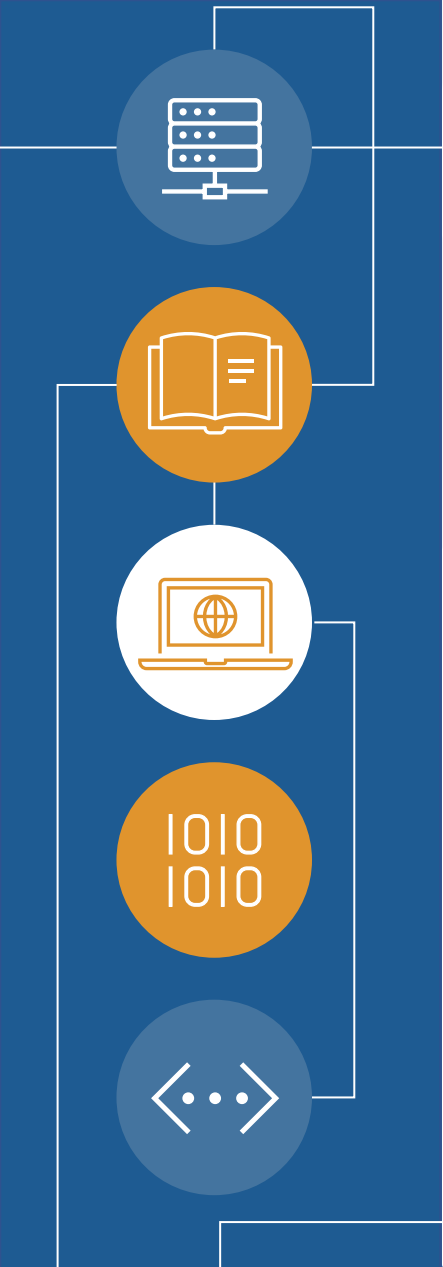
Dengue virus



Data Interoperability: status and challenges

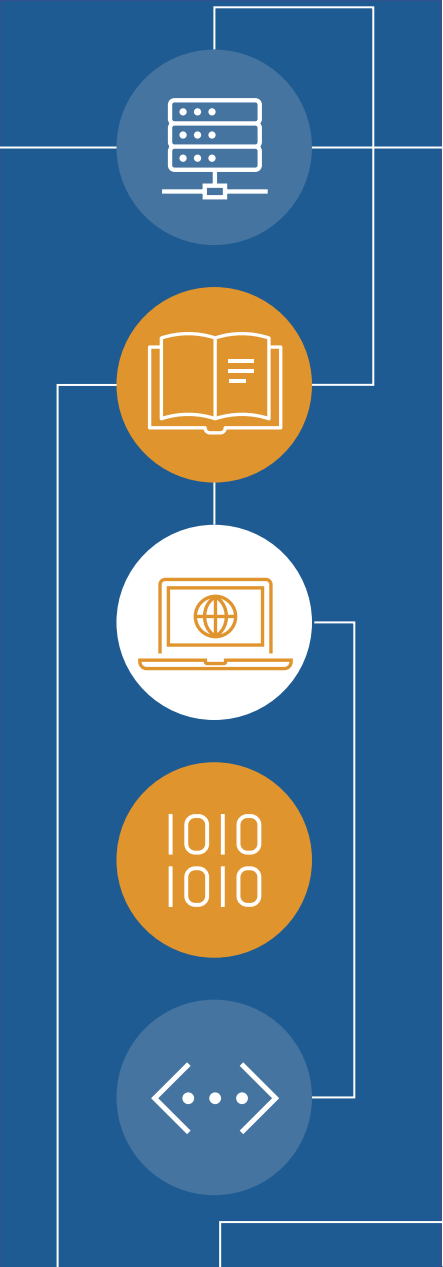
Public access repositories received ~50% of global data generated during COVID-19 pandemic

- Standards
 - Variable
 - Minimal required information (collection burden vs. utility)
- Data Exchange
 - Limited by restrictive repository egress policies
- Lack of oversight/regulation
 - No performance goals for data interoperability

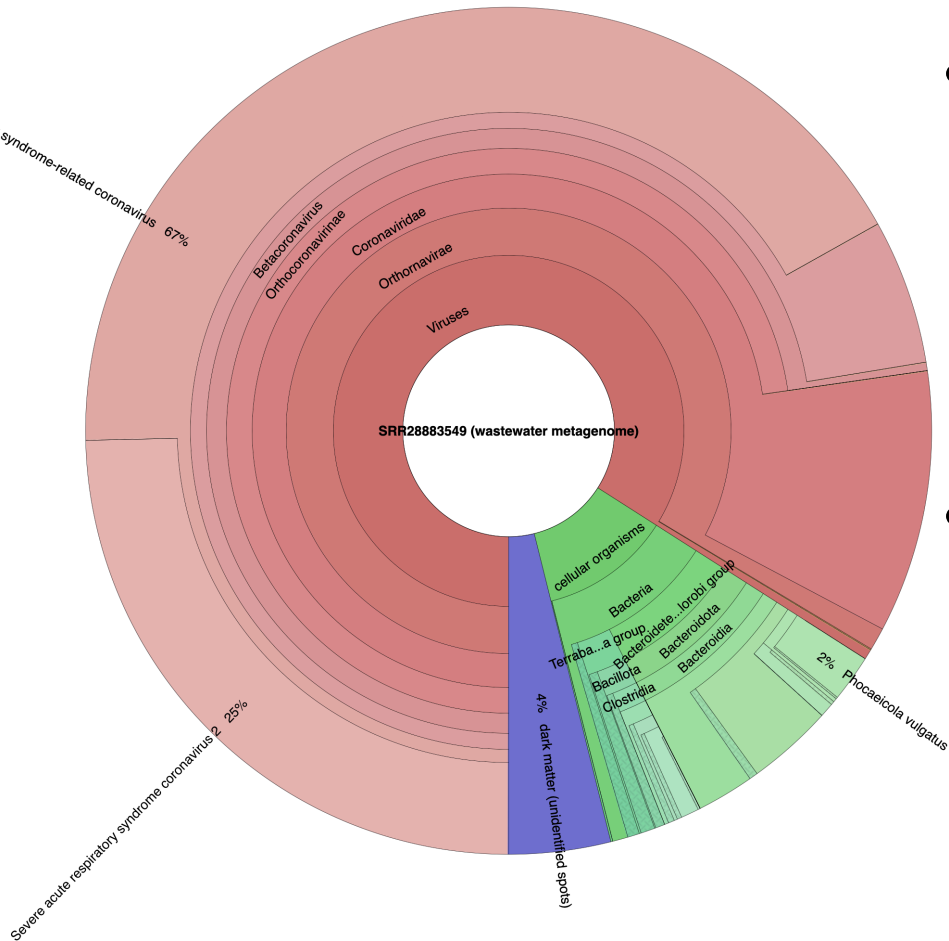


Pathogen Genomics Infrastructure: Status and Challenges

- Distribution
 - Interoperability
 - Latency
 - Scalability
 - Repository access policies
 - Public communications timelines
- Sustainability: current trends
- Potential opportunities for AI



Wastewater and environmental biosurveillance



- All sequences submitted to SRA analyzed by the Sequence Taxonomic Analysis Tool (STAT)
 - Identifies the organismal content of each sequence sample – species, genera, families
 - Results published on web, AWS, and GCP and can be used to monitor viruses and organisms
- Public Wastewater Surveillance Consortium
 - Formed by NCBI to foster interagency collaboration
 - Wastewater and environmental data, metadata, and analysis standards and best practices

Katz KS, et al. *Genome Biol.* 2021 Sep 20;22(1):270. doi: 10.1186/s13059-021-02490-0.

NCBI pathogen genomics lessons learned

- SRA and GenBank are critical to sharing SARS-CoV-2 sequence data and descriptive metadata with the scientific and public health communities
- SRA data critical to data validation and scientific insights
- Standardized data reduction key to making large scale data useful for downstream analyses
- Community tools like the NCBI developed “human sequence scrubber” necessary to support pathogen data
- Partnership-driven analysis best practices reduce burden on individual researchers or groups



1010
1010

1010
1010

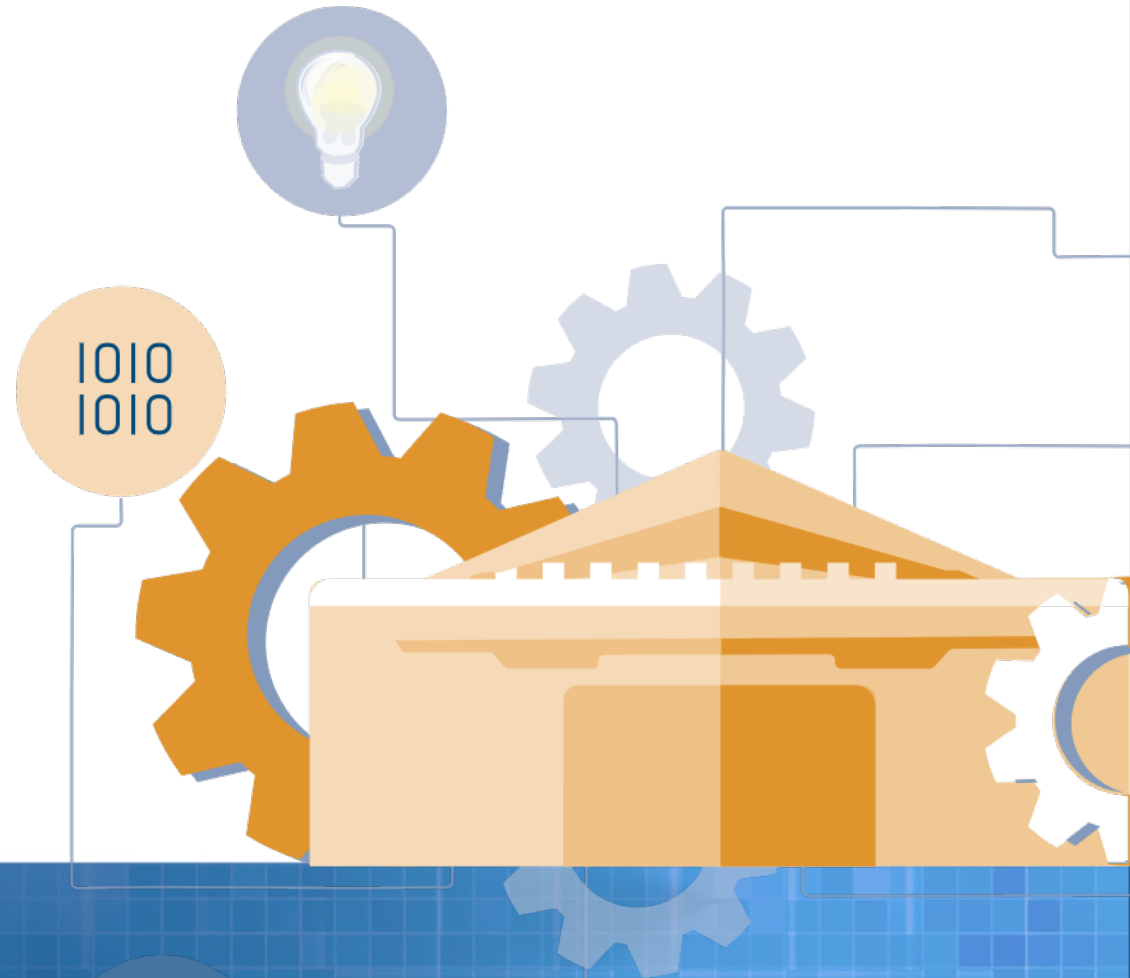
Questions?



National Library of Medicine
National Center for Biotechnology Information



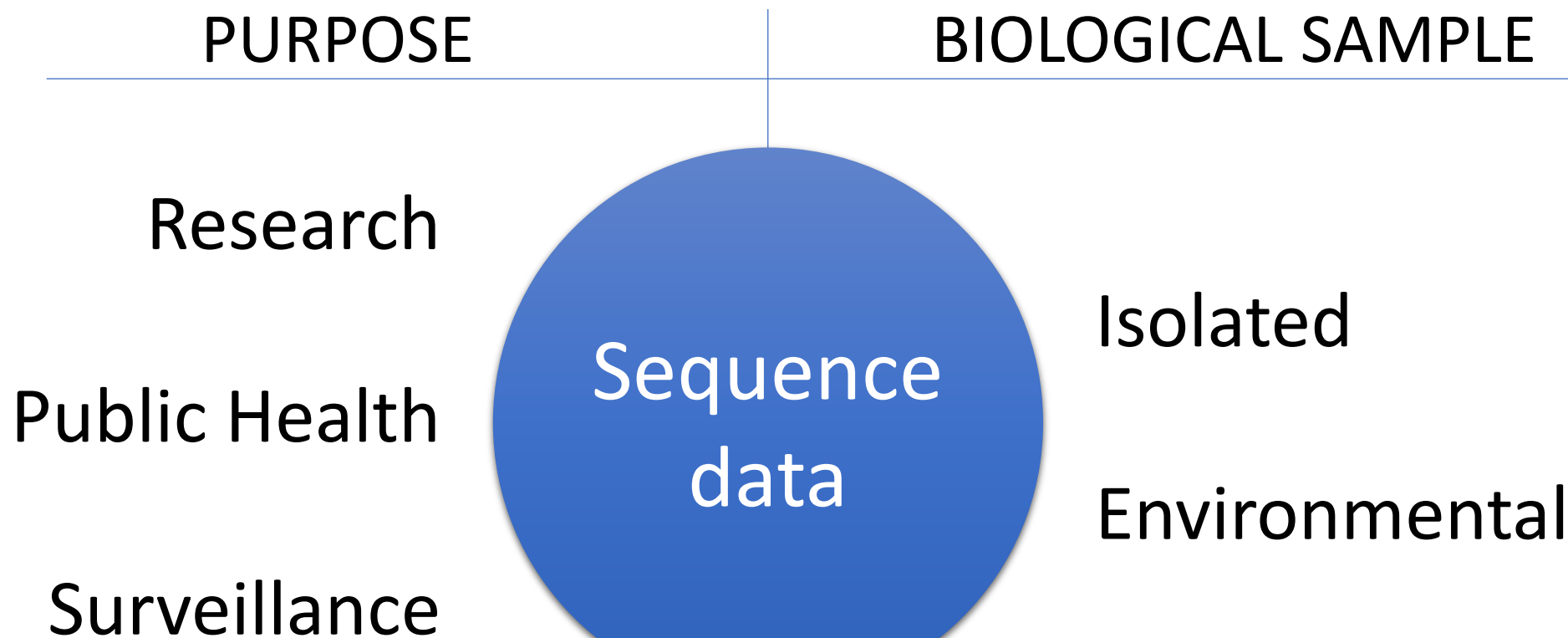
Extra Slides



- NCBI**



Sequence Data Generation



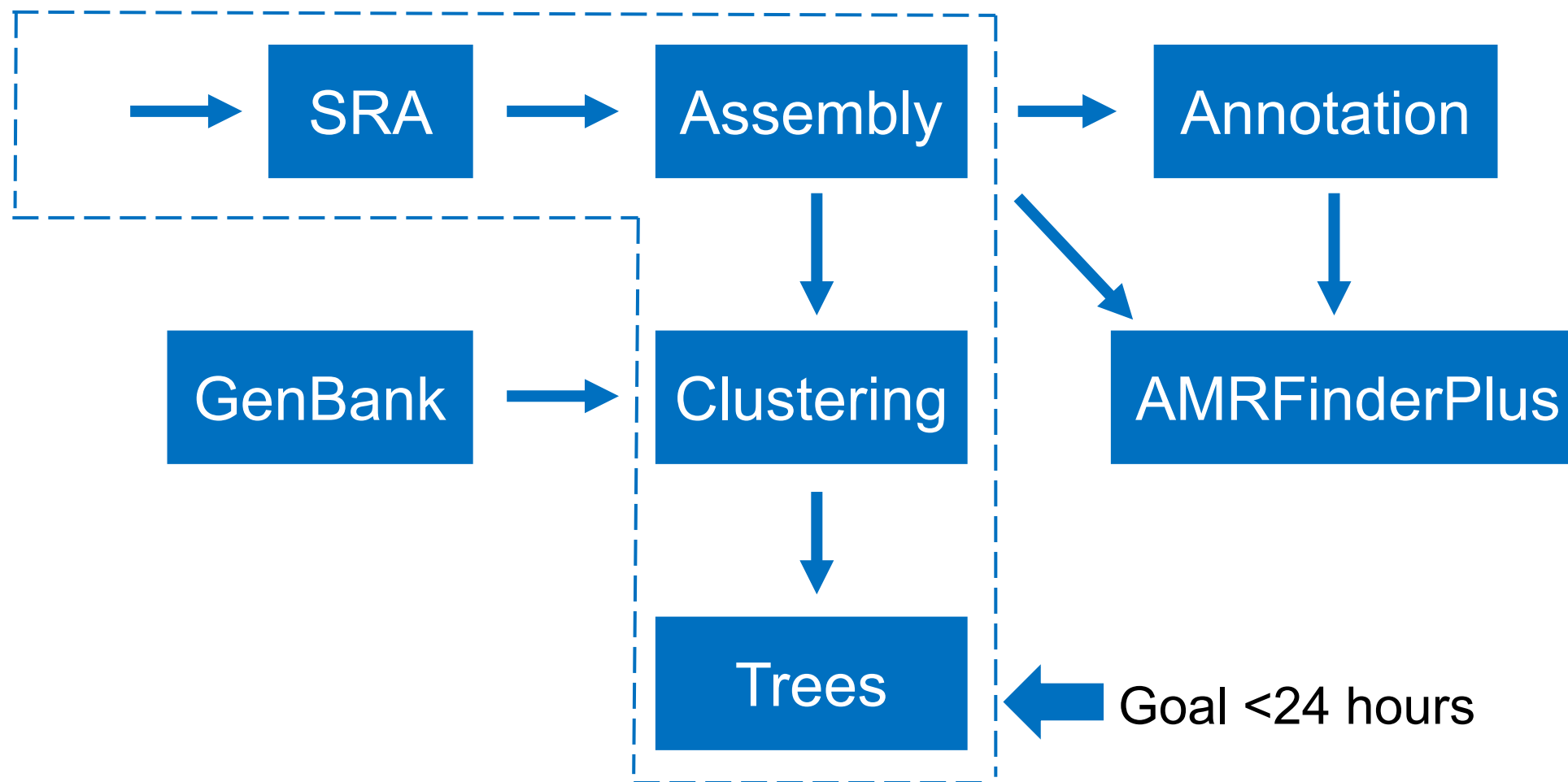
Combating Antibiotic Resistant Bacteria (CARB)

To address the growing threat of antibiotic resistance, the U.S. Government released the National Strategy for CARB in September 2014

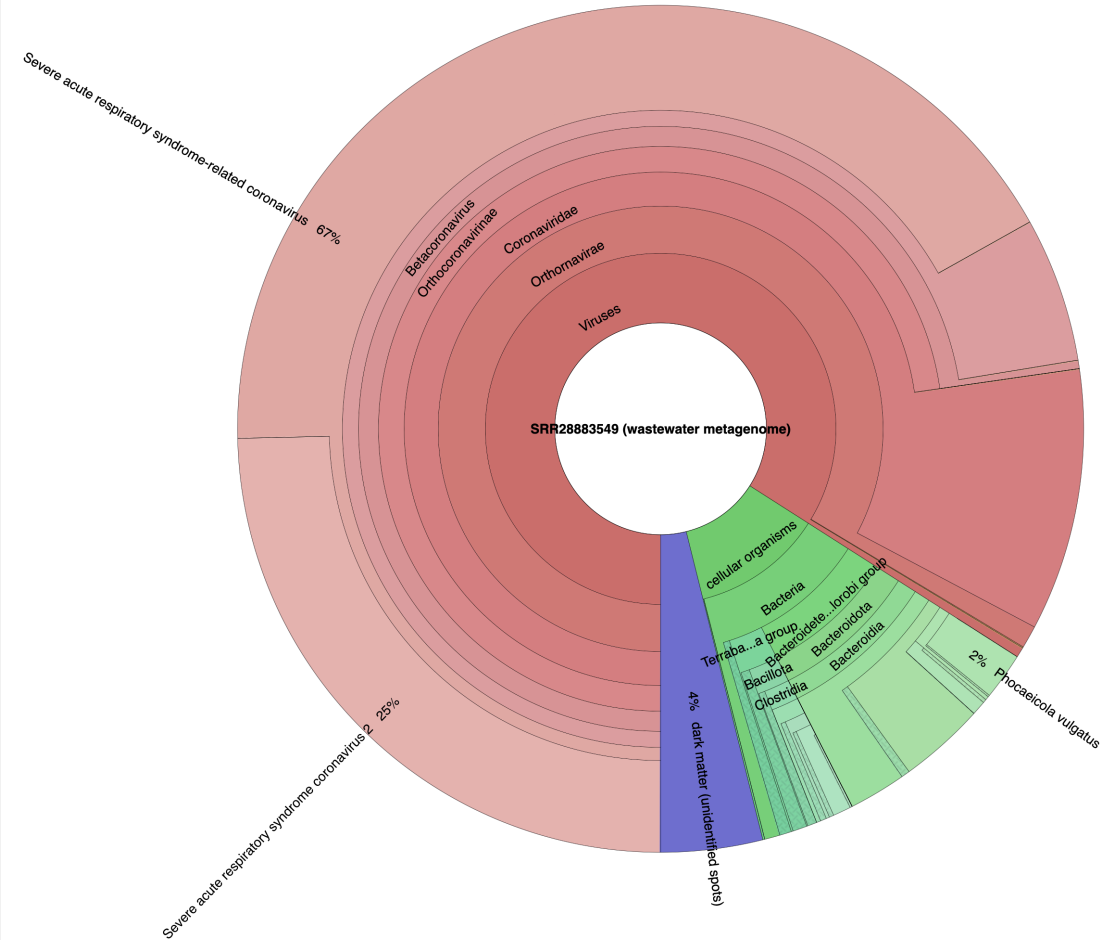
NLM receives support via NIAID for this effort

- Developed the National Database of Antibiotic-Resistant Organisms (NDARO) - a centralized hub for researchers to access AR data to facilitate real-time surveillance of pathogenic organisms.
 - developed and maintains a curated reference database of AR, virulence, and stress response genes
 - developed AMRFinderPlus to identify these genes in bacterial genomes
 - provides web interfaces on the presence of these genes in more than 1.9 million isolates along with associated metadata

Pathogen Detection (PD) Pipeline



Wastewater and environmental biosurveillance

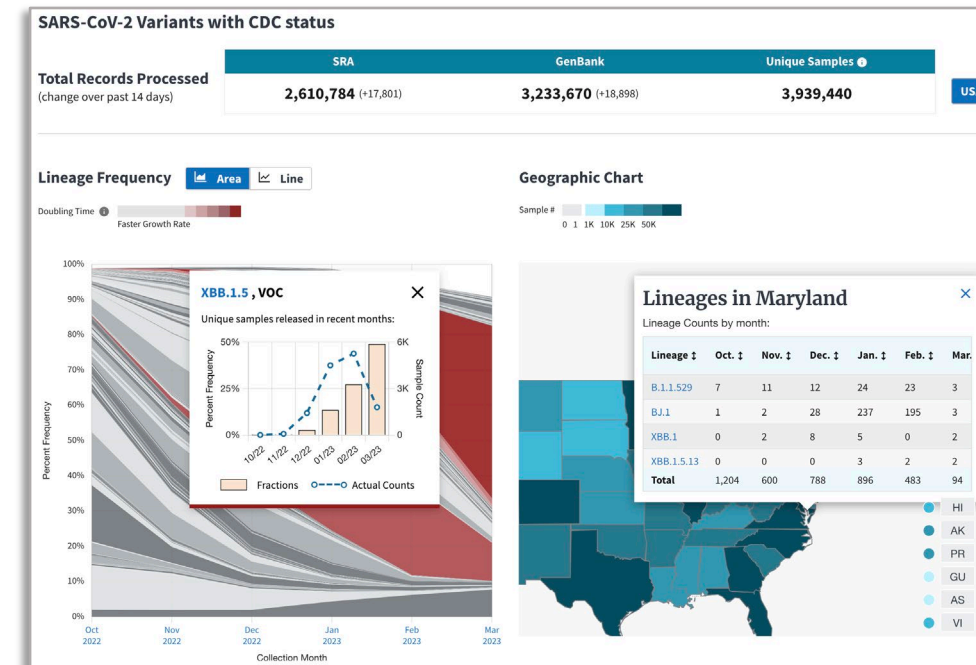


- All sequences submitted to SRA analyzed by the Sequence Taxonomic Analysis Tool (STAT)
 - Identifies the organismal content of each sequence sample – species, genera, families, etc.
 - Results published on web, AWS, and GCP and can be used to monitor viruses and organisms
- Public Wastewater Surveillance Consortium
 - Formed by NCBI to foster interagency collaboration
 - Wastewater and environmental data, metadata, and analysis standards and best practices

Katz KS, et al. *Genome Biol.* 2021 Sep 20;22(1):270. doi: 10.1186/s13059-021-02490-0.

Lessons learned from NCBI pathogen genomics support

- Partnership-driven analysis best practices reduce burden on individual researchers or groups
- SRA and GenBank were critical to sharing SARS-CoV-2 sequence data and descriptive metadata with the scientific and public health communities
- SRA data critical to data quality and validation
 - Mutations, linkage, and recombination
- Standardized data reduction key to making the data useful for downstream analyses
 - 500 TB compressed into in easy-to-use tables
 - 1.4 M pathogen samples compressed into actionable data in short time frame
- Tools like the NCBI developed “human sequence scrubber” necessary to support pathogen data



Selected viral genomic submissions to INSDC in the past year

Virus	GenBank Records			SRA Records		
	No. of Records	No. of Countries	% Missing Location	No. of Records	No. of Countries	% Missing Location
SARS-CoV-2	1,080,335	58	0.0%	428,591	59	0.3%
Influenza A virus	144,282	57	2.0%	6,370	20	4.4%
Influenza B virus	27,448	16	0.6%	368	6	0.0%
Dengue virus	5,416	72	0.0%	1,610	40	0.2%
African swine fever virus	709	35	0.4%	63	10	20.6%
Measles virus	1,104	25	0.0%	54	3	5.6%
Ebolaviruses	207	2	0.0%	14	2	0.0%
Marburg virus	1	1	0.0%	8	2	0.0%
Rift Valley fever virus	101	8	1.0%	68	3	0.0%
Zika virus	141	11	6.0%	674	4	15.4%
HIV-1	51,391	49	12.2%	1,590	11	9.1%