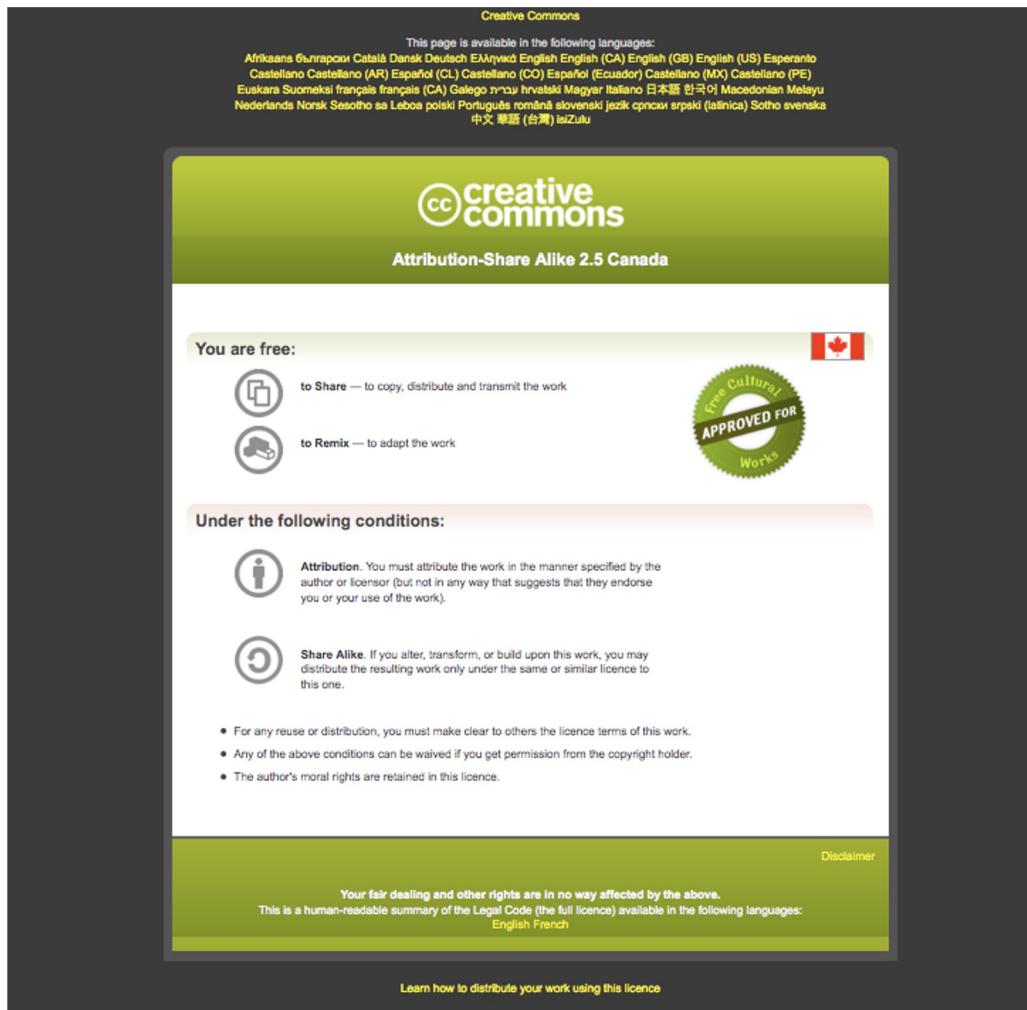




Canadian Bioinformatics Workshops

www.bioinformatics.ca

bioinformaticsdotca.github.io



Module 1

Introduction to Genomic Epidemiology



William Hsiao

Infectious Disease Genomic Epidemiology

April 18-21, 2023



CENTRE FOR
INFECTIOUS DISEASE
GENOMICS AND
ONE HEALTH



SIMON FRASER
UNIVERSITY

Course Overview

- Module 1: **Introduction** – definitions, backgrounds, purpose and benefits of genomic epidemiology
- Module 2: **Phylogenetic Analysis**
- Module 3: **Data Curation and Data Sharing**
- Module 4: **Viral Pathogen Genomic Analysis (SNV)**
- Module 5: **Bacterial Pathogen Genomic Analysis (gene level)**
- Module 6: **Antimicrobial Resistant Gene Analysis**
- Module 7: **Phyldynamics and Transmission Dynamics**
- Module 8: **Emerging Pathogen Detection and Identification**
- Module 9: **Mobile Genetic Elements and Environmental Microbiome**
- Keynote: **Samira Mubareka (TBC)**

General Learning Objectives

- By the end of this lecture, you will:
 - Understand how genomic epidemiology can improve clinical and public health microbiology
 - Process genomic sequence data using a variety of bioinformatics tools for bacterial and viral genomes and metagenomes
 - Interpret genomic data in epidemiological context and understand the importance of data standardization and sharing
 - Perform several types of genomic epidemiology analyses
 - Recognize the limitations and challenges of genomic epidemiology (a rapidly evolving field)

Learning Objectives of Module 1

- Understand why infectious disease research is important
- Be familiar with some examples of genomic epidemiology studies
- Be familiar with high-throughput sequencing and its application to clinical and public health microbiology
- Be familiar with sequence data processing
- Understand the challenges associated with sharing genomic epidemiology data

Global Flight Paths

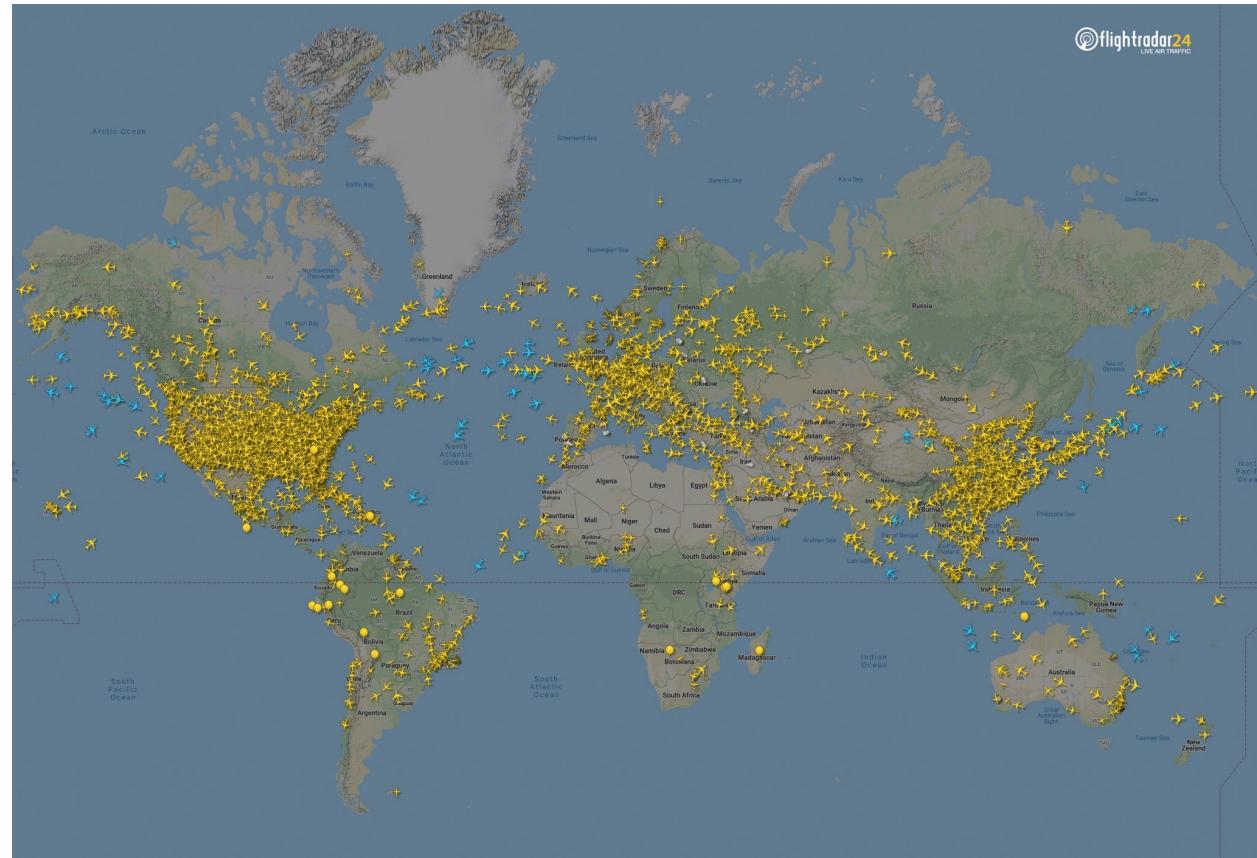


<http://openflights.org/data.html>

March 7 2020

v.s.

April 7 2020



[Then and now: visualizing COVID-19's impact on air traffic | Flightradar24 Blog](#)

OPEN

Meta-genomic analysis of toilet waste from long distance flights; a step towards global surveillance of infectious diseases and antimicrobial resistance



Received: 17 December 2014

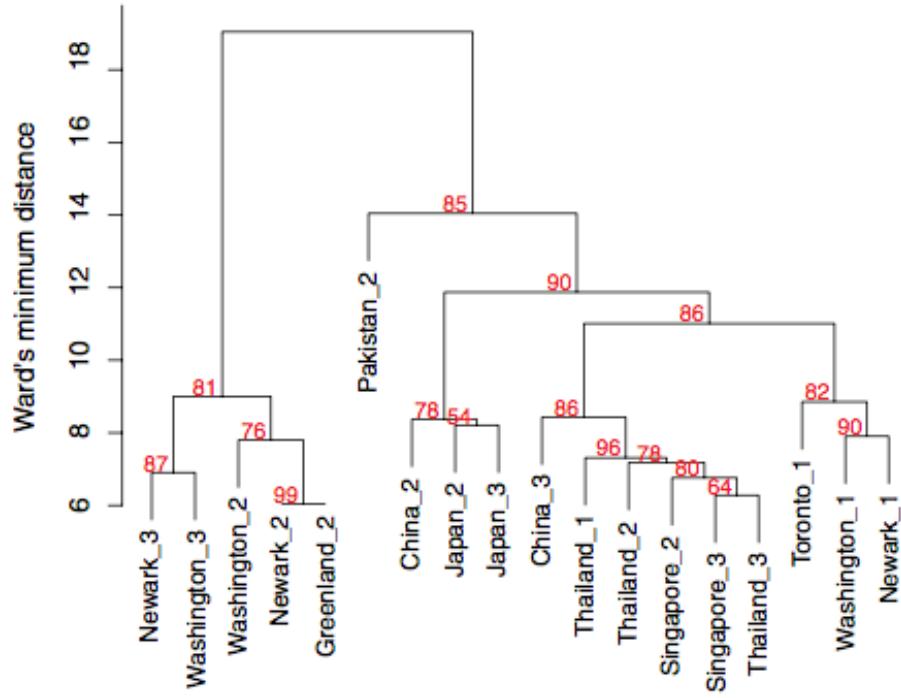
Accepted: 17 April 2015

Published: 10 July 2015

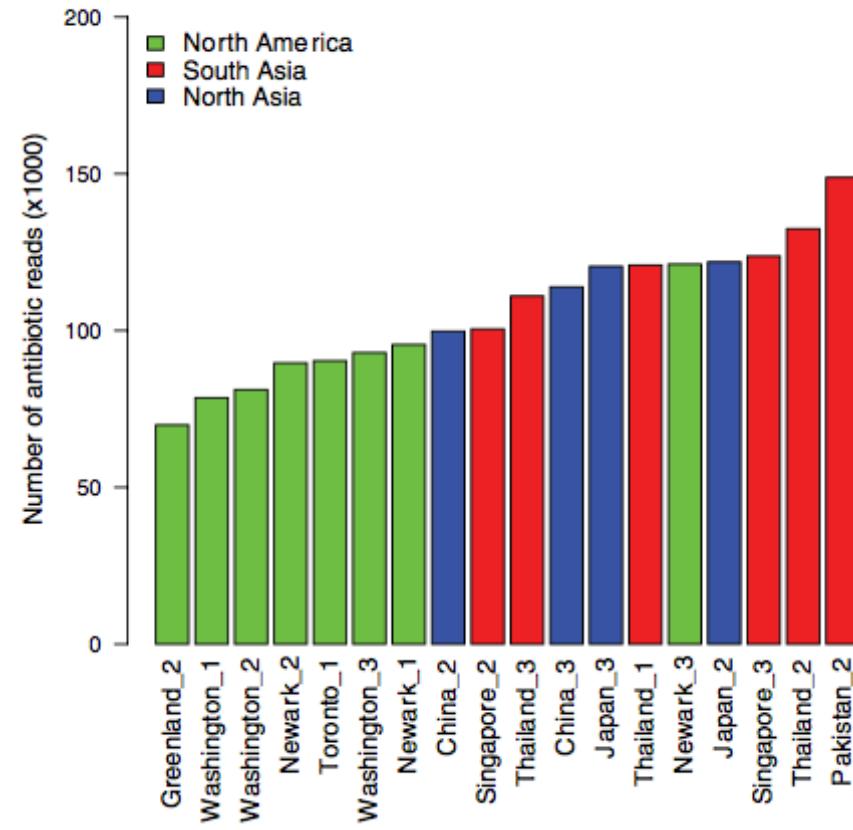
Thomas Nordahl Petersen¹, Simon Rasmussen¹, Henrik Hasman², Christian Carøe¹, Jacob Bælum¹, Anna Charlotte Schultz², Lasse Bergmark², Christina A. Svendsen², Ole Lund¹, Thomas Sicheritz-Pontén¹ & Frank M. Aarestrup²

- 18 flights from 3 continents
- Onboard human wastes (400L per flight!) extracted for DNA sequencing
- Samples clustered based on microbiome profile
- Antimicrobial resistance genes identified (~0.06% of the reads)

Petersen, T. N. *et al.* Meta-genomic analysis of toilet waste from long distance flights; a step towards global surveillance of infectious diseases and antimicrobial resistance. *Sci Rep* 1–9 (2015). doi:10.1038/srep11444



Clustering by geographic origins



Higher proportions of antibiotic resistance genes found in flights from South Asia

Petersen, T. N. et al. Meta-genomic analysis of toilet waste from long distance flights; a step towards global surveillance of infectious diseases and antimicrobial resistance. *Sci Rep* 1–9 (2015). doi:10.1038/srep11444

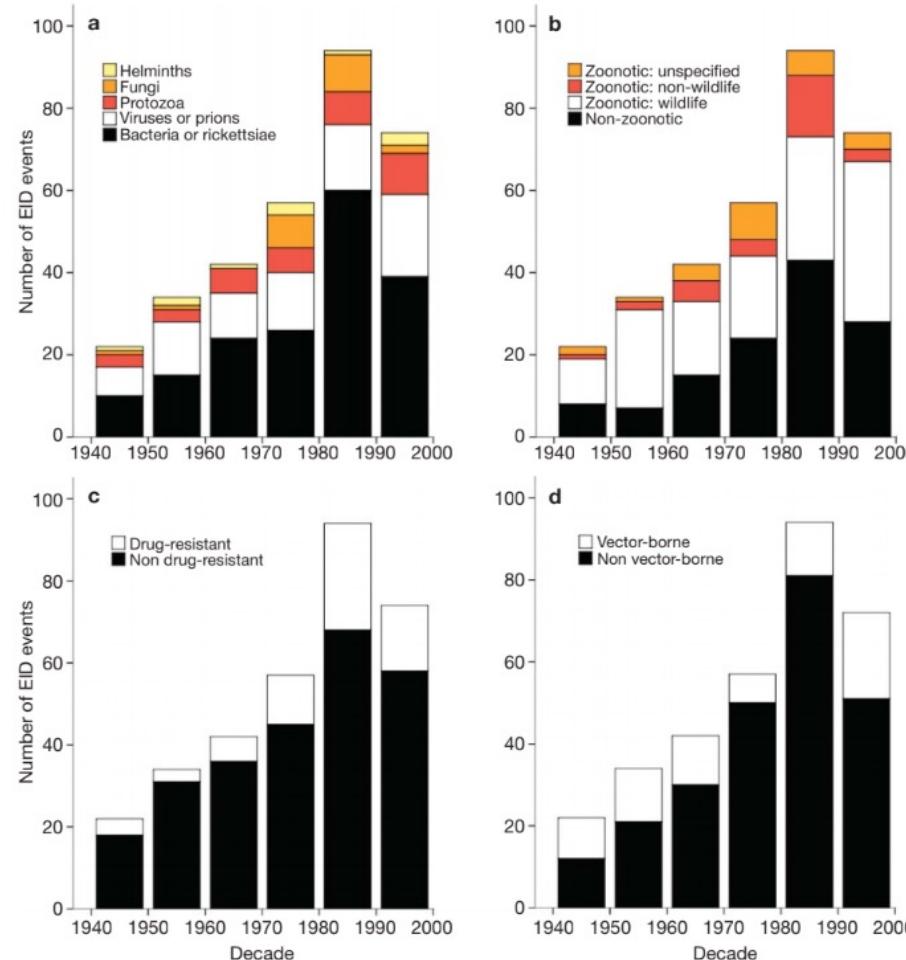
Increase in Emerging Infectious Diseases

EID events: detection of newly evolved strains of pathogens in human

Dominated by zoonotic diseases

Identified global hotspots for EID events

Identified risk factors for EIDs: tropical rain forest, population density, climate, mammal species richness, etc.

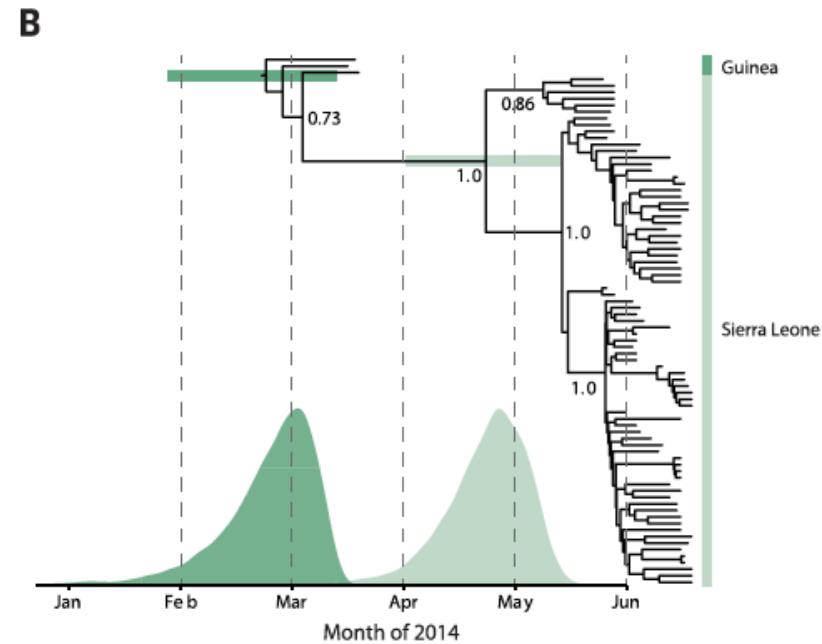


Jones. et al 2008. "Global Trends in Emerging Infectious Diseases." *Nature* 451 (7181): 990–93.

Allen. 2017. "Global Hotspots and Correlates of Emerging Zoonotic Diseases." *Nature Communications* 8 (1): 1124.

Genomic Sequence Analysis of Ebola

- Outbreak in West Africa from Dec 2013 - May 2016 resulted in 28,616 reported cases and 11,310 deaths
- Global impact on travel
- Genomic epidemiological analysis of 99 early isolates revealed:
 - A single human exposure to natural reservoir (likely bats)
 - Outbreak sustained by **human to human transmission** (e.g. funeral practice and lack of proper quarantine facilities)
 - Transmission from Guinea to Sierra Leone likely to be from a single event but two distinct lineages



- Sequence data corroborated with epidemiological narrative critical to unravel the complex situation and institute effective policies and interventions... Still there were significant delays

<http://www.who.int/csr/disease/ebola/en/>

Gire et al. 2014. *Science* 345 (6202): 1369–72.

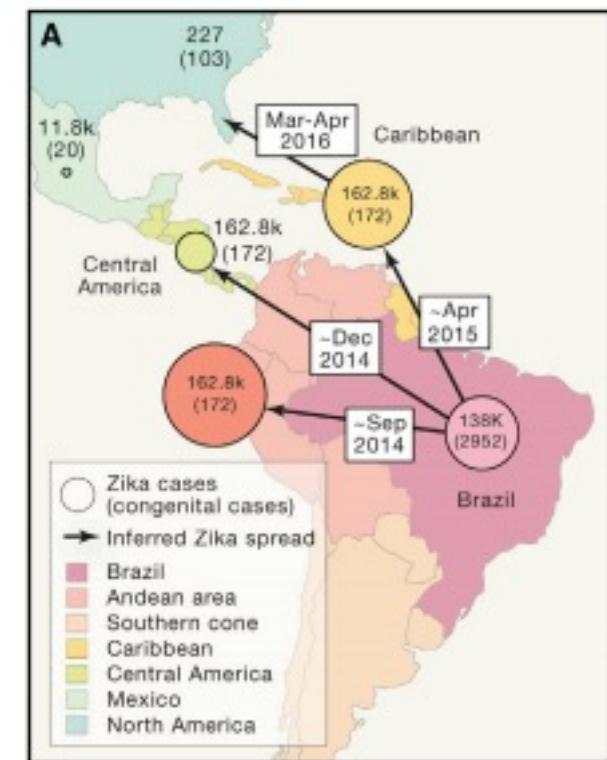
Fast forward to 2018

- Another Ebola outbreak in Democratic Republic of Congo
 - DRC has experienced Ebola outbreaks before so more aware of the symptoms and how to deal with patients
 - Fast mobilization of resources
 - WHO/World Bank has now standing emergency fund to deal with public health emergencies
 - Stockpile of experimental Ebola vaccine (developed at NML, Canada) trialed in previous outbreak available to people at most risk (health workers, contacts of Ebola cases)
 - Experimental drugs available
 - Outside health care workers able to be in position faster

<https://www.cnn.com/2018/05/29/health/ebola-outbreak-2018-response-explainer/index.html>

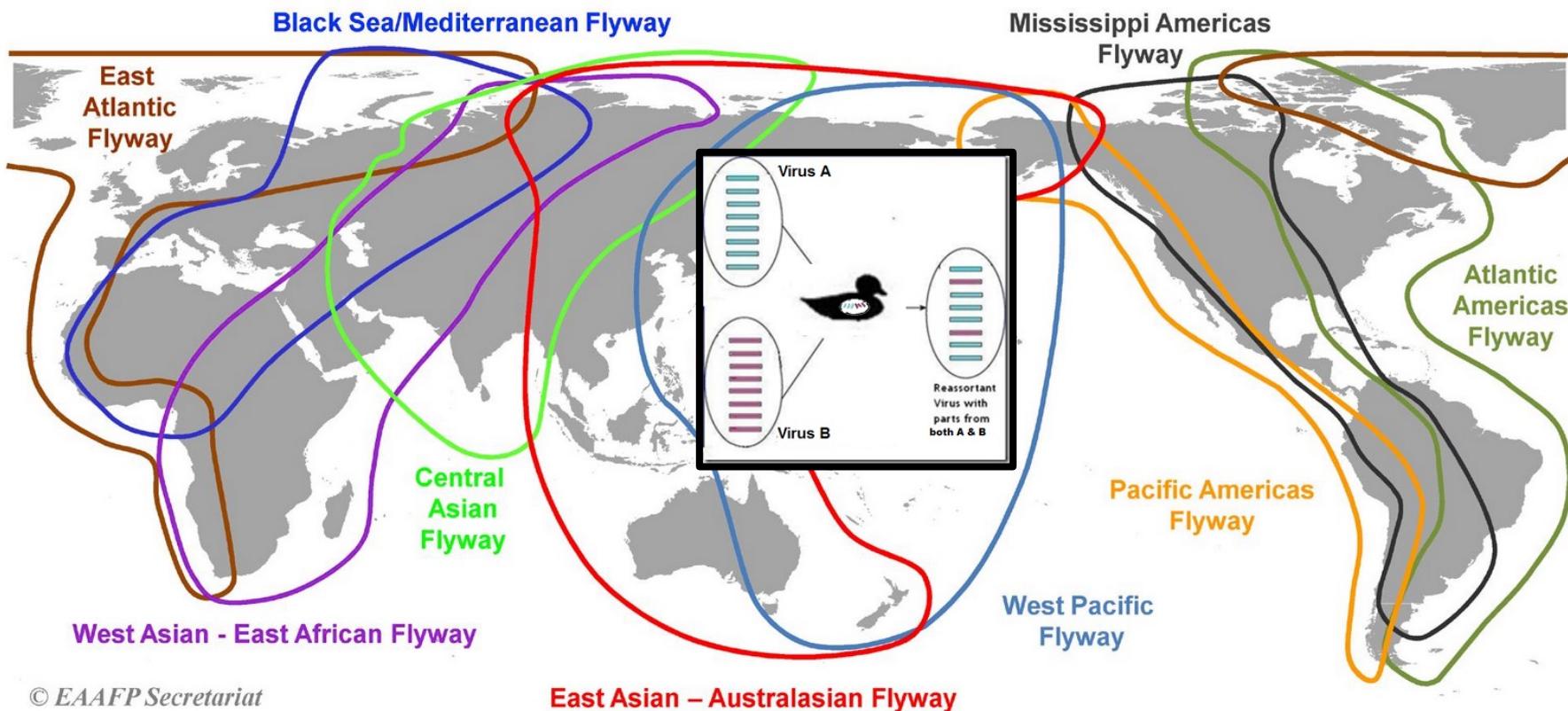
Genomic Insights into Zika Outbreak

- Endemic in Africa and Asia; mild symptoms; self-eliminating
- Caused an outbreak in Americas in 2015–2016; naïve population; high number of people infected and resulted in a few thousand microcephaly cases
- Due to symptoms overlapping with other viruses (e.g. dengue virus), syndromic surveillance can be un-reliable
 - Serological or genomic/molecular tests needed for confirmation
- Phylogenetic reconstruction using genomic data highlighted unsuspected circulation of zika in the population prior to outbreak and help to reconstruct the transmission route

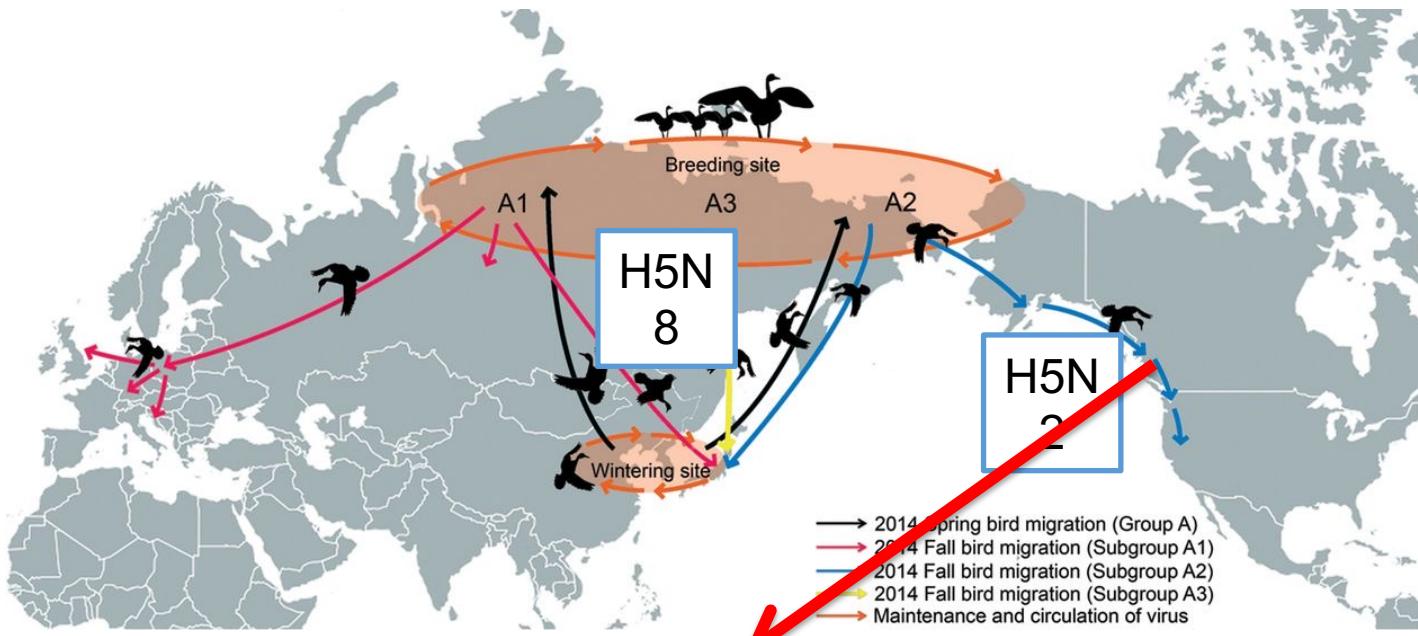


Grubaugh. 2018. "Genomic Insights into Zika Virus Emergence and Spread." *Cell* 172 (6): 1160–62.

Genomic Sequence Analysis of Avian Flu



- Influenza A viruses caused 4 human pandemics in the past 100 years – new hypervirulent strains arise from mixing of human and animal flu viruses
- Birds are natural reservoir to influenza viruses



Current Approaches for Surveillance

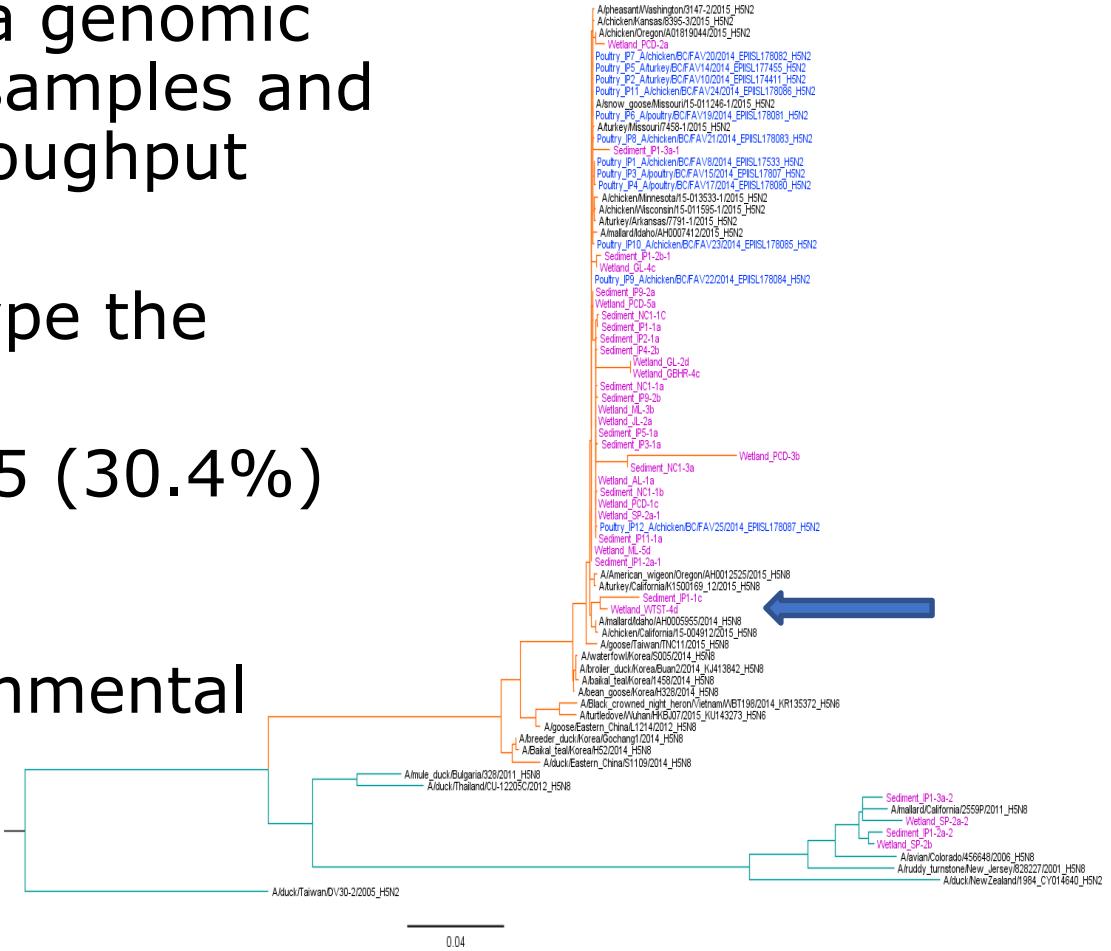
- Passive
 - Testing waterfowl dead from other causes
- Active
 - Capture and testing of live birds
 - Hunter-kill birds
- Overall positivity rate < 1%

Failed to detect the presence of AI in waterfowl in advance of the 2015 outbreak.



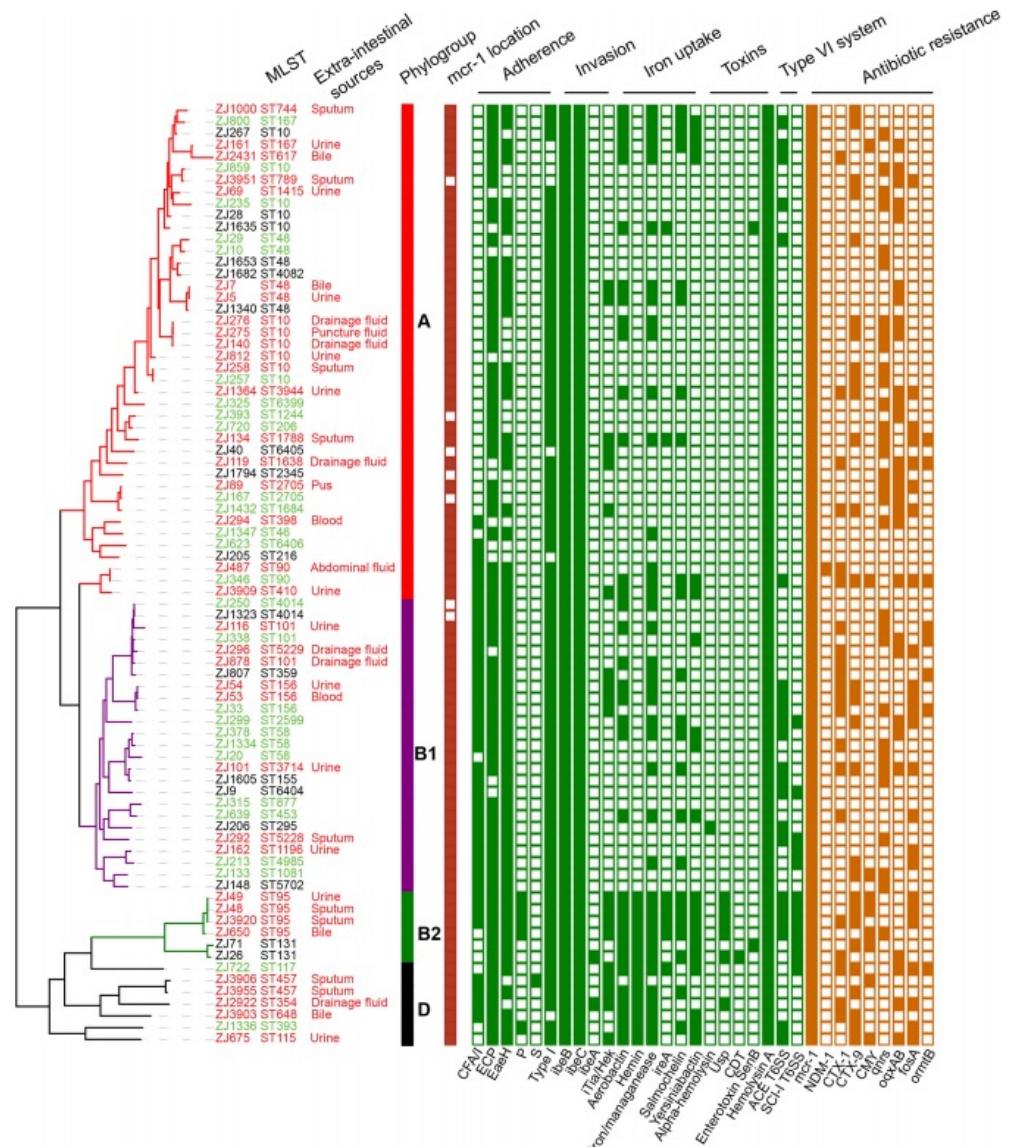
Environmental Genomic Based Surveillance

- Isolate and enrich for influenza genomic RNAs from wetland sediment samples and sequence them using high-throughput sequencing (Illumina)
- Bioinformatic analysis to subtype the sequences
- Overall positivity rate: 105/345 (30.4%)
 - Wetland samples: 72/300 (24.0%)
 - Farm samples: 33/45 (73.3%)
- Able to identify HPAI in environmental samples



AMR and Mobile Elements

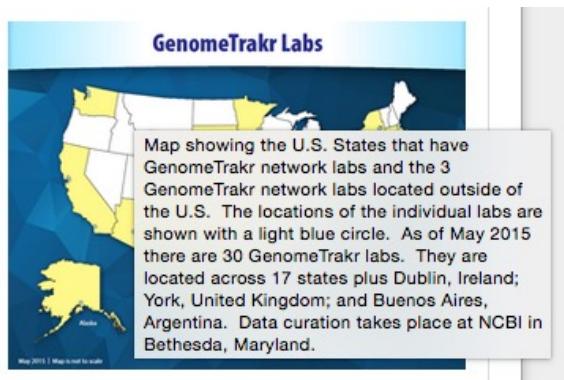
- Many antimicrobial resistant (AMR) genes can move around hosts (mobile elements)
 - Detection of the gene by PCR or identification of the organism alone are insufficient to understand the transmission of these AMR genes



Shen et al. 2018. "Heterogeneous and Flexible Transmission of Mcr-1 in Hospital-Associated Escherichia Coli." <https://doi.org/10.1128/mBio.00943-18>.

Whole Genome Sequencing of Foodborne Pathogens

- UK Public Health England sequence all the *Salmonella* isolates submitted to PH Lab since April 2014
- US FDA and CDC (supported by National Center for Biotechnology Information) created a distributed network of labs to utilize WGS for pathogen identification with a centralized analysis platform (GenomeTrakr)



GOV.UK

Blog
Public health matters

Organisations: Public Health England

A revolution in genomic sequencing

Christine McCartney, 20 January 2014 — Microbiology services, Protecting the country's health

Microbiology is the study of microscopic organisms such as viruses and bacteria. In the past as microbiologists we would have relied on growing...

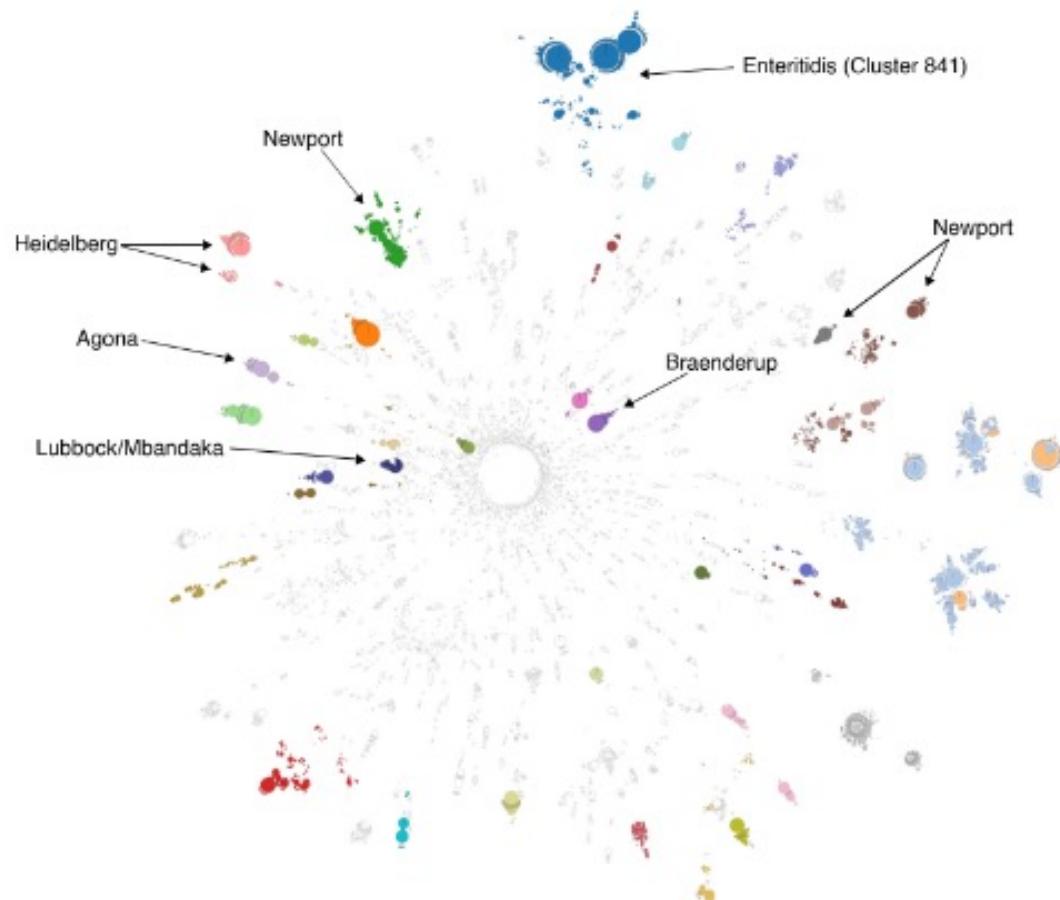
Search blog

Public health matters

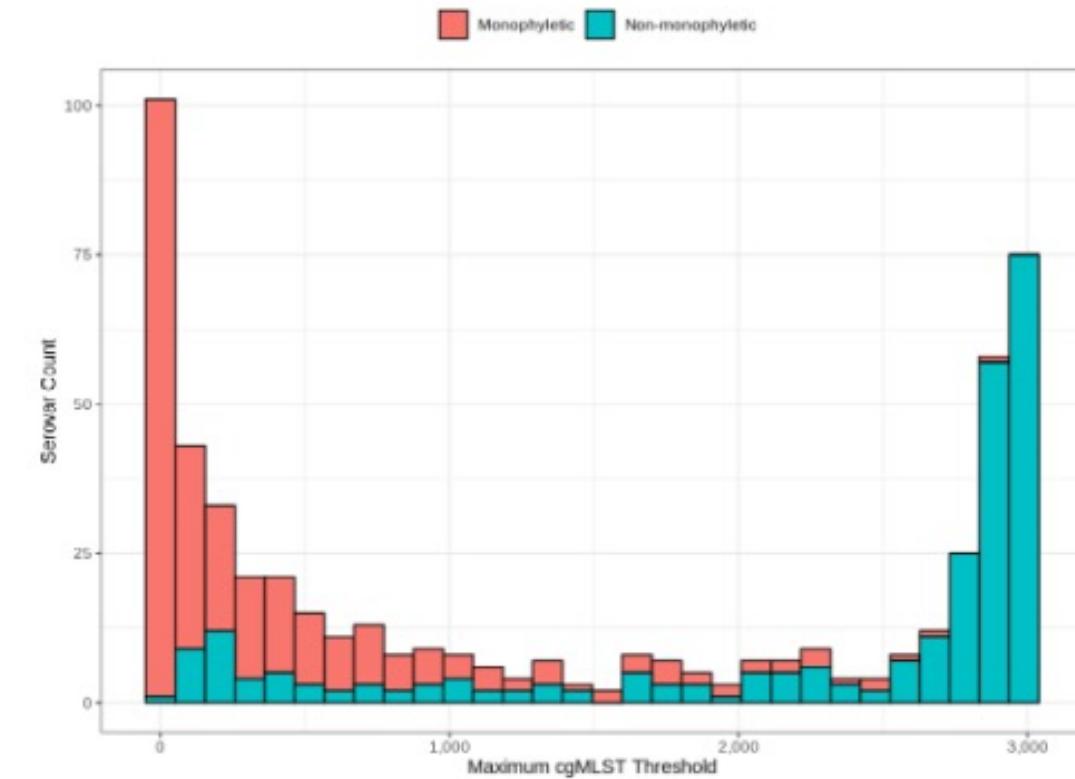
The official blog of Public Health England, providing expert insight on the organisation's work and all aspects of public health. More

<https://publichealthmatters.blog.gov.uk/2014/01/20/innovations-in-genomic-sequencing/>
<http://www.fda.gov/Food/FoodScienceResearch/WholeGenomeSequencingProgramWGS/ucm363134.htm>

Large-scale comparative genomics to refine the organization of the global *Salmonella enterica* population structure

Chao Chun Liu^{1,2} and William W. L. Hsiao^{1,2,3*}

- Cluster-Serovar**
- 841-enteritidis [33854]
 - 70-typhimurium [1989]
 - 786-infantis [9010]
 - 70-i 1,4,[5],12:i- [779]
 - 801-newport [7613]
 - 624-kentucky [4787]
 - 276-javiana [4368]
 - 775-heidelberg [4134]
 - 53-braenderup [3298]
 - 689-agona [3120]
 - 63-newport [2822]
 - 72-saintpaul [2788]
 - 31-thompson [2420]
 - 114-anatum [2302]
 - 84-newport [2161]
 - 99-muenchen [2112]
 - 204-montevideo [1971]
 - 235-oranienburg [1844]
 - 257-schweizergrund [1803]
 - 838-dublin [1719]
 - 666-mbandaka [1683]
 - 655-senftenberg [1625]
 - 84-hadar [1410]
 - 48-stanley [1323]
 - 103-muenchen [1295]
 - 742-derby [1262]
 - 691-weltevreden [1145]
 - 224-readng [1120]
 - 520-paratyphi a [1044]
 - 482-mississippi [1029]
 - 202-montevideo [1018]
 - 584-kentucky [952]
 - 93-vinchow [908]
 - 837-beria [855]
 - 213-panama [847]
 - Others [42358]



Genomics has been a hero of the COVID-19 pandemic.

Cite as: X. Deng *et al.*, *Science*
10.1126/science.abb9263 (2020).

A SARS-CoV-2 vaccine candidate would likely match all currently circulating variants

Bethany Dearlove, Eric Lewitus, Hongjun Bai, Yifan Li, Daniel B. Reeves, M. Gordon Joyce, Paul T. Scott, Mihret F. Amare, Sandhya Vasan, Nelson L. Michael, Kayvon Modjarrad, and Morgane Rolland

PNAS September 22, 2020 117 (38) 23652-23662; first published August 31, 2020;
<https://doi.org/10.1073/pnas.2008281117>

Check for updates

The proximal origin of SARS-CoV-2

Kristian G. Andersen, Andrew Rambaut, W. Ian Lipkin, Edward C. Holmes & Robert F. Garry

Nature Medicine 26, 450–452(2020) | Cite this article

5.03m Accesses | 706 Citations | 35003 Altmetric | Metrics

To the Editor – Since the first reports of novel pneumonia (COVID-19) in Wuhan, Hubei province, China^{1,2}, there has been considerable discussion on the origin of the causative virus, SARS-CoV-2³ (also referred to as HCoV-19)⁴. Infections with SARS-CoV-2 are now widespread, and as of 11 March 2020, 121,564 cases have been confirmed in more than 110 countries, with 4,373 deaths⁵.

SARS-CoV-2 is the seventh coronavirus known to infect humans; SARS-CoV, MERS-CoV and SARS-CoV-2 can cause severe disease, whereas HKU1, NL63, OC43 and 229E are associated with mild symptoms⁶. Here we review what can be deduced about the origin of SARS-CoV-2 from comparative analysis of genomic data. We offer a perspective on the notable features of the SARS-CoV-2 genome and discuss scenarios by which they could have arisen. Our analyses clearly show that SARS-CoV-2 is not a laboratory construct or a purposefully manipulated virus.

Genomic surveillance reveals multiple introductions of SARS-CoV-2 into Northern California

Xianding Deng^{1,2*}, Wei Gu^{1,2*}, Scot Federman^{1,2*}, Louis du Plessis^{3*}, Oliver G. Pybus², Nuno Faria³, Candace Wang^{1,2}, Guixia Yu^{1,2}, Brian Bushnell⁴, Chao-Yang Pan⁵, Hugo Guevara⁵, Alicia Sotomayor-Gonzalez^{1,2}, Kelsey Zorn⁶, Allan Lopez⁷, Venice Servellita⁸, Elaine Hsu¹, Steve Miller¹, Trevor Bedford^{1,8}, Alexander L. Greninger^{7,9}, Pavitra Roychoudhury^{7,9}, Lea M. Starita^{8,10}, Michael Famulare¹, Helen Y. Chu^{8,11}, Jay Shendure^{8,9,12}, Keith R. Jerome^{7,9}, Catie Anderson¹⁴, Karthik Gangavarapu¹⁴, Mark Zeller¹⁴, Emily Spencer¹⁴, Kristian G. Andersen¹⁴, Duncan MacCannell¹⁵, Clinton R. Paden¹, Yan Li¹⁵, Jing Zhang¹⁵, Suxiang Tong¹⁵, Gregory Armstrong¹⁵, Scott Morrow¹⁶, Matthew Willis¹⁷, Bela T. Matyas¹⁸, Sundari Mase¹⁹, Olivia Kasirye²⁰, Maggie Park²¹, Godfred Masinde²², Curtis Chan²², Alexander T. Yu², Shua J. Chai^{5,15}, Elsa Villarino²³, Brandon Bonin²³, Debra A. Wadford²³, Charles Y. Chiu^{1,2,24†}

Comment on this paper

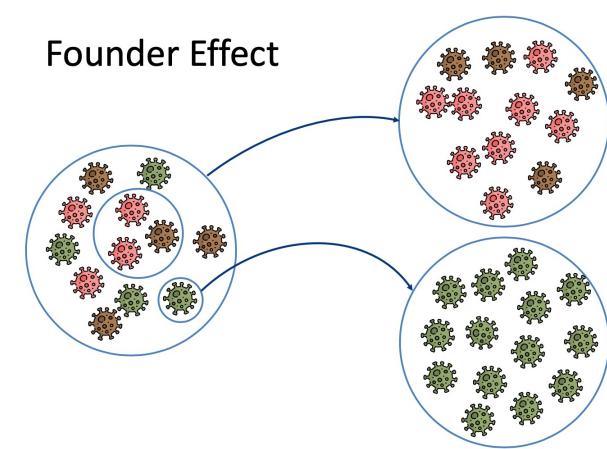
Large scale sequencing of SARS-CoV-2 genomes from one region allows detailed epidemiology and enables local outbreak management

Andrew J Page, Alison E Mather, Thanh Le Viet, Emma J Meader, Nabil-Fareed J Alikhan, Gemma L Kay, Leonardo de Oliveira Martins, Alp Aydin, David J Baker, Alexander J. Trotter, Steven Rudder, Ana P Tedim, Anastasia Kolyva, Rachael Stanley, Maria Diaz, Will Potter, Claire Stuart, Lizzie Meadows, Andrew Bell, Ana Victoria Gutierrez, Nicholas M Thomson, Evelien M Adriaenssens, Tracey Swinler, Rachel Aj Gilroy, Luke Griffith, Dheeraj K Sethi, Rose K Davidson, Robert A Kingsley, Luke Bedford, Lindsay J Coupland, Ian G Charles, Ngozi Elumogo, John Wain, Reenesh Prakash, Mark A Webber, SJ Louise Smith, Meera Chand, Samir Dervisevic, Justin O'Grady, The COVID-19 Genomics UK (COG-UK) consortium

doi: <https://doi.org/10.1101/2020.09.28.20201475>

Why sequence SARS-CoV-2 genomes

- Transmission tracking at the regional, provincial, national and international scales
- Cluster investigations (i.e. genomic epidemiology)
- Evolving viral characteristics that might impact
 - detection methods (PCR, serology)
 - clinical outcomes (strain severity) and transmission
 - The most reliable way to detect variants of concern
 - effectiveness of healthcare measures, treatments and vaccines (ID regions of sequence frequently changing – or not changing at all – informing drug target/vaccine development and drug/vaccine continued effectiveness)
- Many national and continental efforts to sequence the virus in close-to-real-time (COG-UK, SPHERES, Australia, Pan-Africa, etc) – however some of these efforts are now being sun setted



Canadian COVID-19 Genomics Network

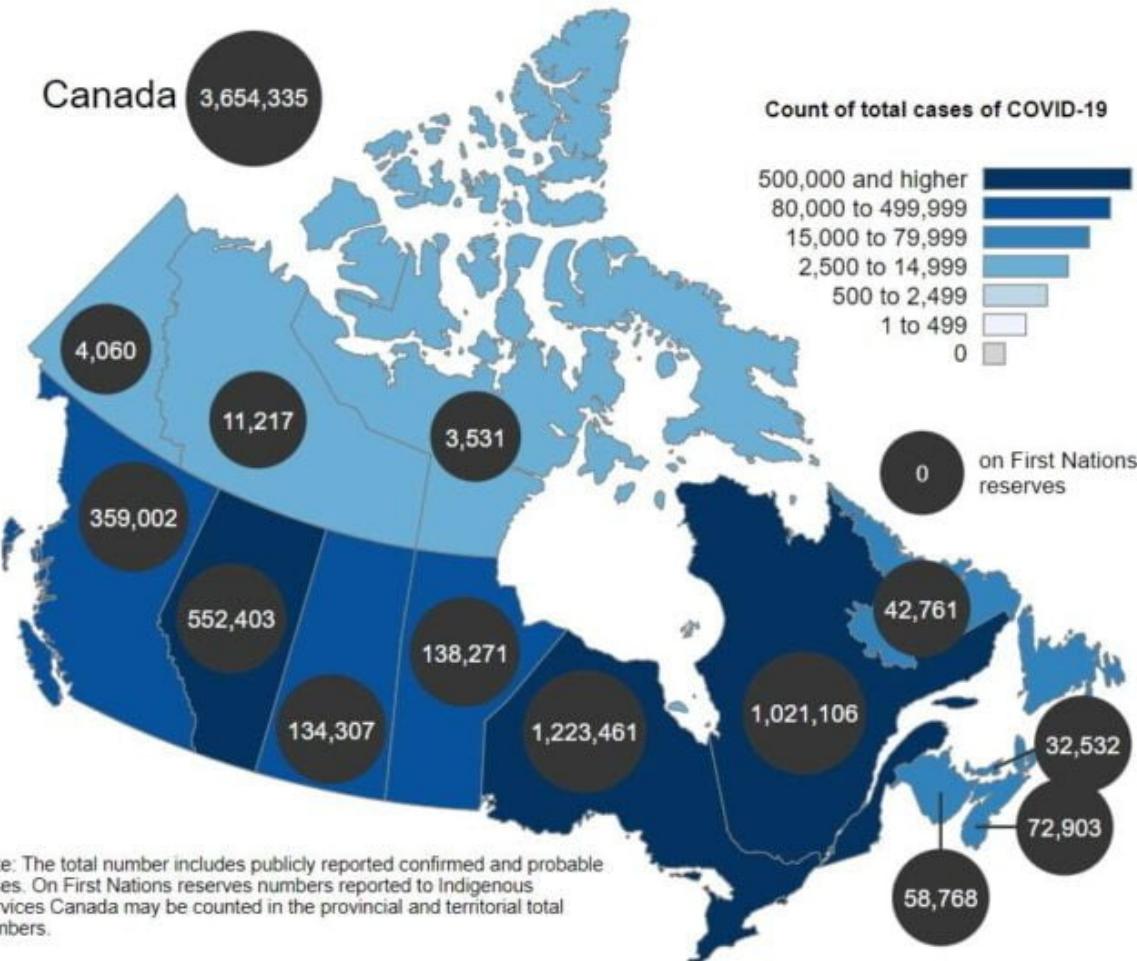
- Established in March 2020 with an initial \$40 Million Canadian Federal Government Investment
- \$20 million viral genomic sequencing and genomic capacity building in public health
- \$20 million human host genome sequencing from infected individuals
- Consortium of national and provincial public health laboratories, hospitals, research institutes, large-scale genome sequencing centers, industry, and coordinating centers.
- **Goals**
 - Coordinate and fund SARS-CoV-2 and human host genome sequencing efforts (up to 150,000 viral; 10,000 human)
 - Integrate **sequence data** and harmonize associated **clinical/epidemiological data (metadata)** – led by my group
 - Facilitate **data sharing** nationally and internationally
 - Capacity building, including for **future outbreak/pandemic preparedness**

<https://www.genomecanada.ca/en/cancogen>



The Canadian COVID-19 Genomics Network

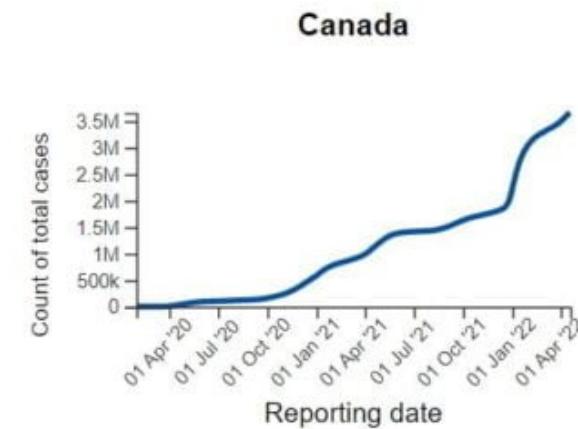
- 150K viral genomes, 10K human genomes



Count of total cases of COVID-19



The count of total cases of COVID-19 in Canada was 3,654,335 as of April 19, 2022.



To date: >500,000 viral genomes were sequenced in Canada

VirusSeq Data Portal – publicly accessible viral genomics data from Canada

The screenshot shows the VirusSeq Data Portal homepage. At the top, there is a navigation bar with the VirusSeq logo, followed by links for "Explore VirusSeq Data", "Analysis Tools", "About", and "Data Releases". Below the navigation bar, a large red Canadian flag icon is positioned next to the text "Canadian VirusSeq Data Portal". A descriptive paragraph explains the project's goal: "The goal of the CanCOGeN VirusSeq project was to sequence up to 150,000 viral samples from Canadians testing positive for COVID-19. The VirusSeq Data Portal is an open-source and open-access data portal for all Canadian SARS-CoV-2 sequences and associated non-personal contextual data. It harmonizes, validates and automates submission to international databases." Below this text, four key statistics are displayed: "471,185 Files", "471,185 Viral Genomes", "12 Studies", and "14 GB". At the bottom, two teal-colored buttons are visible: "Explore the Data" and "Download the Data".

Canadian VirusSeq Data Portal

The goal of the CanCOGeN VirusSeq project was to sequence up to 150,000 viral samples from Canadians testing positive for COVID-19. The VirusSeq Data Portal is an open-source and open-access data portal for all Canadian SARS-CoV-2 sequences and associated non-personal contextual data. It harmonizes, validates and automates submission to international databases.

471,185 Files **471,185** Viral Genomes **12** Studies **14** GB

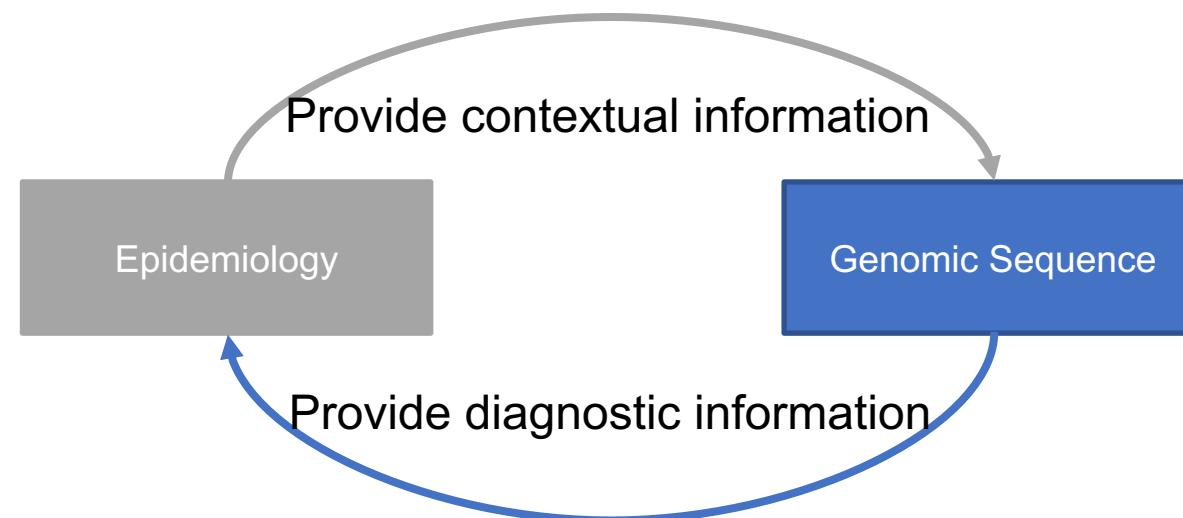
[Explore the Data](#) [Download the Data](#)

(McGill, OICR, SFU – leveraging existing software platform solutions, the portal was created within 3 months of being funded)

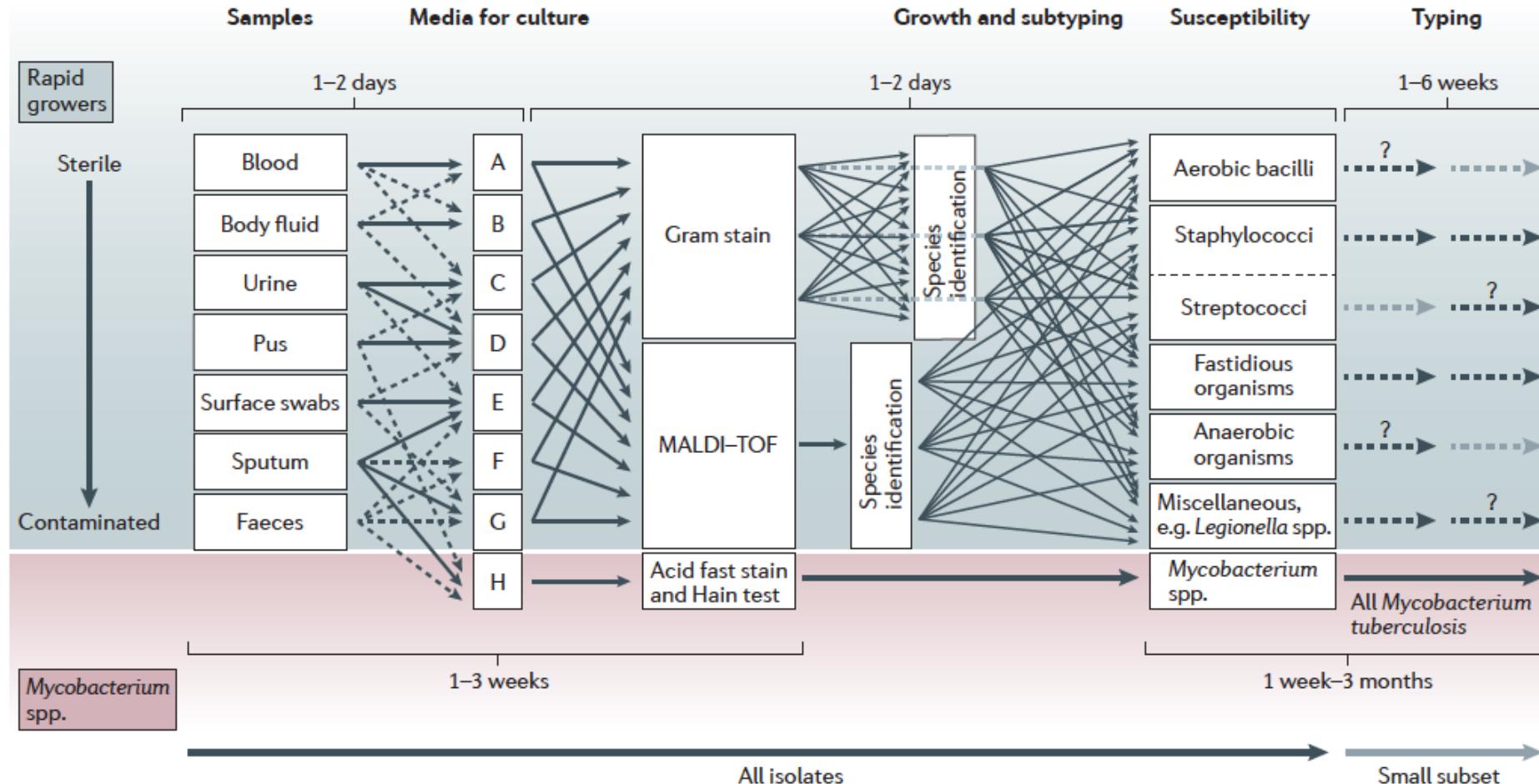
bioinformatics.ca

Genomic Epidemiology

- Def: Combine **whole genome sequencing** data from **pathogens** with **epidemiological investigations** to track spread of an infectious disease



Traditional Clinical Microbiology Laboratory



Didelot et al. 2012. doi:10.1038/nrg3226.

Traditional Methods of Characterizing Foodborne Pathogens in a Public Health Laboratory

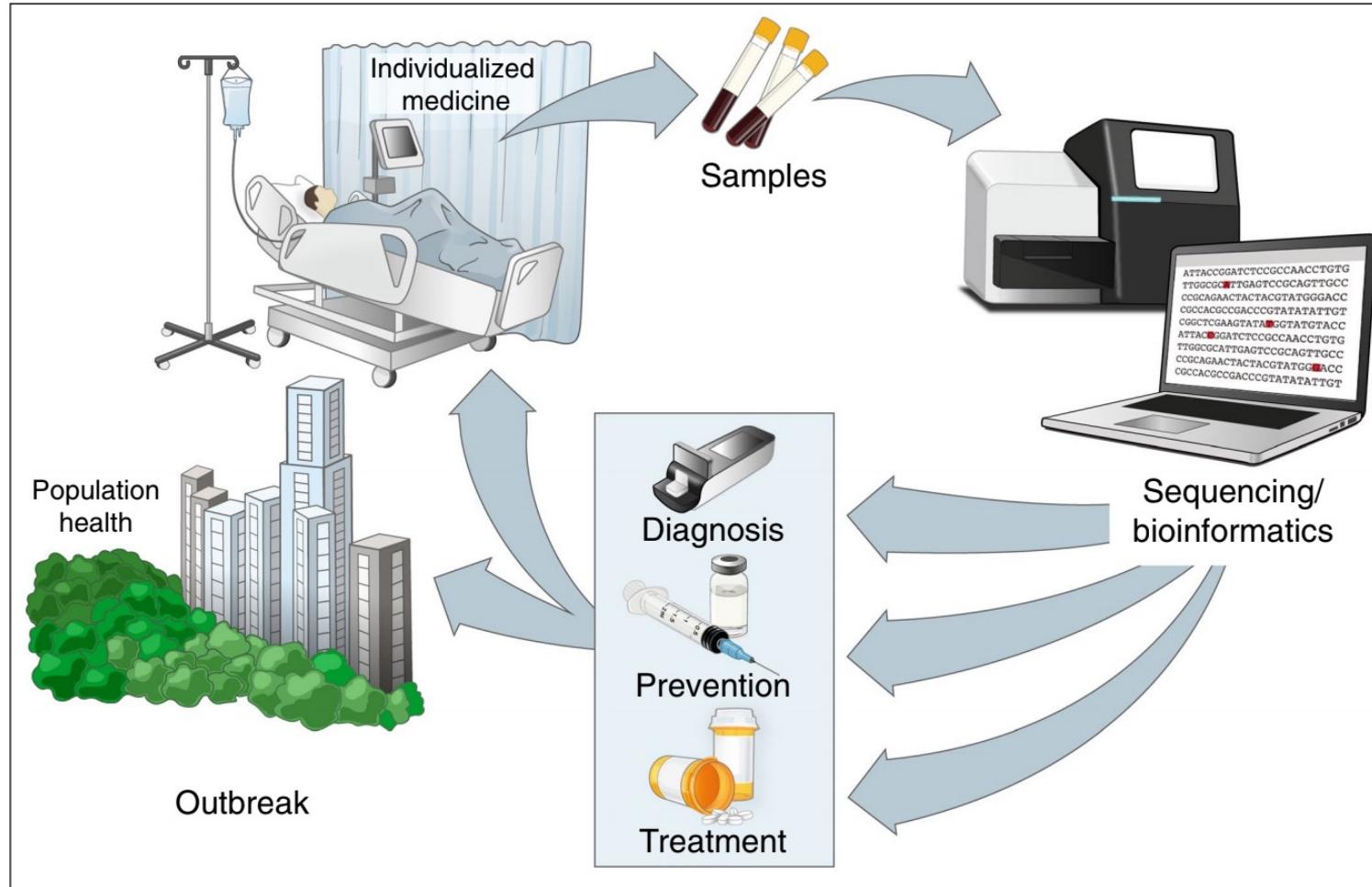
- Growth characteristics
- Phenotypic panels
- Agglutination reactions
- Enzyme immuno assays (EIAs)
- PCR
- DNA arrays (hybridization)
- Sanger sequencing of marker genes
- DNA restriction
- Electrophoresis (PFGE, capillary)

Each pathogen is characterized by methods that are specific to that pathogen **in multiple workflows** (separate workflows for each pathogen)

TAT: 5 min – weeks (months)

Source: Rebecca Lindsey

Whole Genome Sequencing Based Workflow



Ladner et al. 2019 Nature Med

Benefits and Challenges

- Simplified Workflow
- Faster turn-around-time (for some applications)
- Cost-saving by reducing the number of platforms/instruments
- Sequencing is becoming commoditized
- Results (sequences) more **comparable** and **sharable** than other test data
- Value-added analysis (e.g. pathogen evolution, AMR prediction, transmission dynamic modeling etc.)
- Results harder to process and interpret
- Computational Resource Requirement higher (no IT support?)
- Rapid changing technologies
- Per sample cost still higher? Batching required
- Other benefits/challenges?

High Throughput Sequencing (HTS)

- HTS = next-gen sequencing and third-gen sequencing platforms
- Sequence data have many clinical and PH lab utilities
 - Diagnostic – strain level identification, virulence gene ID, AMR gene ID
 - Surveillance –gene-by-gene typing, single-nucleotide-variant (SNV) typing, serotyping, copy-number variants (VNTR, MLVA)
 - Outbreak detection and investigations
 - Source Tracking
- These genomic data can be useful for downstream research use (e.g. improve our understanding of pathogen evolution)



Illumina MiSeq (NGS)



Oxford MinION (3rd Gen)

Sequencing technologies



Sanger



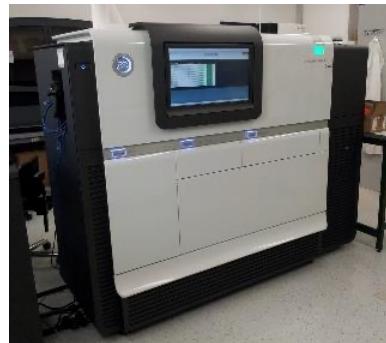
Gene Studio
(Ion Torrent)



Roche 454



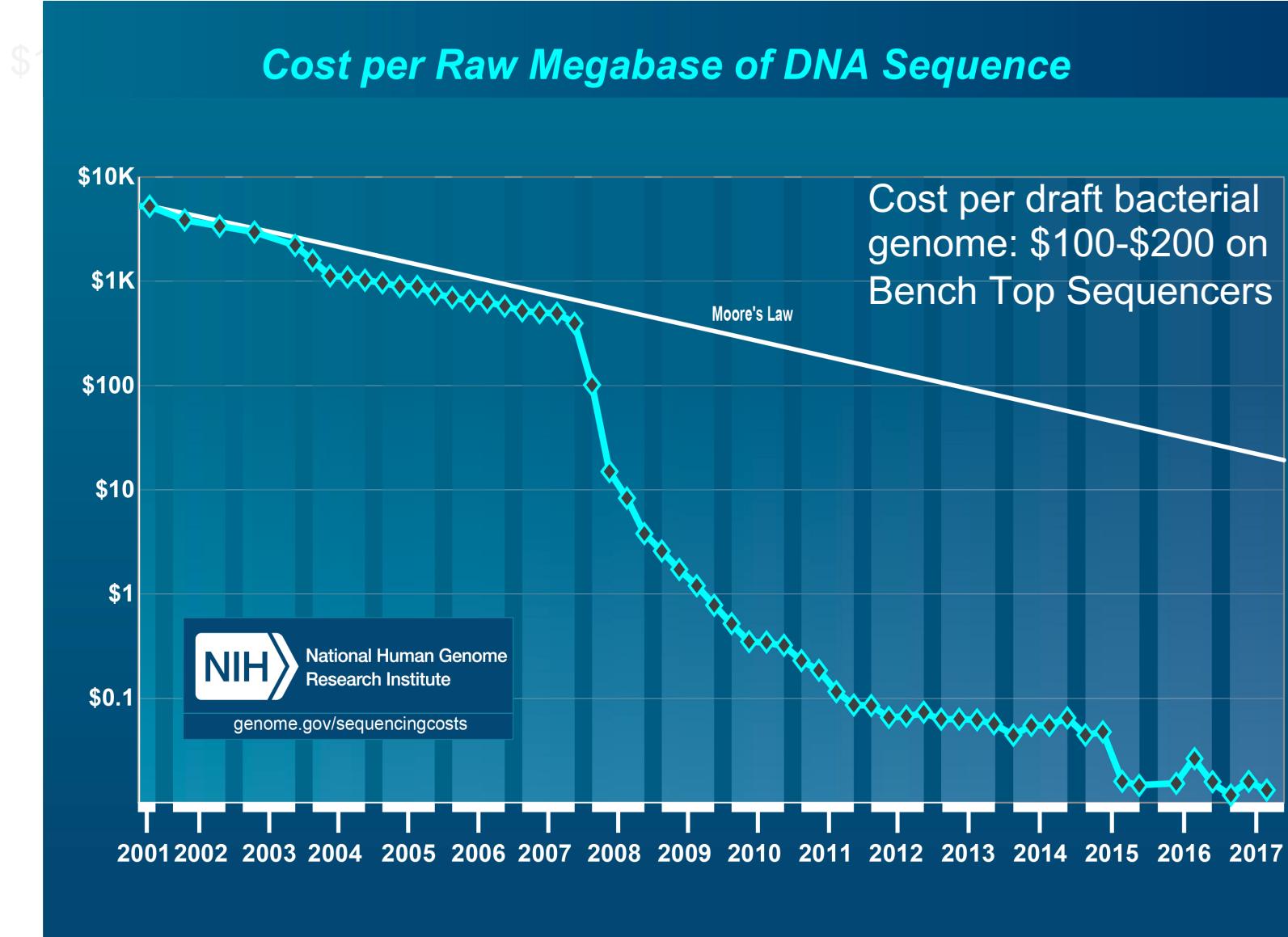
Illumina *Seq



Pacific Biosciences



Nanopore



Wetterstrand KA. DNA Sequencing Costs: Data from the NHGRI Genome Sequencing Program (GSP)
Available at: www.genome.gov/sequencingcostsdata.

Short vs. Long Read Sequencing

Sequencing technology (specifications)	Instrument cost	Library prep cost	Cost per lane/cell	Cost per million reads	Cost per billion bases	Runtime (h)	Throughput (k reads at maximum runtime)	Avg read length (maximum)	Read accuracy (%)	Error profile
PacBio Sequel I (maximum 8 cells)	300,000	150–400	1,000	2,000	100	10–20	500	20,000 (50,000)	85–88/99.9 (CCS)	Homopolymers
PacBio Sequel II (maximum 16 cells)	650,000	75–400	2,000	500	17	10–30	4,000	30,000 (100,000)	88–90/99.9 (CCS)	Homopolymers
ONT Flongle	1,300	40–90	80	400	13	0.1–12	200	30,000 (60,000)	96–99	Homopolymers
ONT MinION	900	90–130	600	600	12	0.1–48	1,000	50,000 (2.3M)	96–99	Homopolymers
ONT GridION (maximum 5 FCs)	45,000	90–130	600	600	12	48	1,000	50,000 (2.3M)	96–99	Homopolymers
ONT PromethION (maximum 48 FCs)	176,000	90–600	1,400	250	8	72	6,000	30,000 (330,000)	96–99	Homopolymers
SLR (example of NovaSeq 2 × 150 PE, 30× coverage, for 5 kb)	900,000	30–50	5,000	1,250	140	44	4,000	9,000 (12,000)	100	Unequal coverage
Illumina NovaSeq S4 (2 × 150 PE; maximum 4 lanes)	900,000	50–100	5,000	2.5	9	44	2,000,000	250 (290)	99.9	Low-quality ends
Illumina MiSeq (2 × 300 PE; maximum 2 lanes)	100,000	50–100	2,000	100	175	56	20,000	550 (590)	99.9	Low-quality ends
ION G5-S5 Prime, P550 chip (maximum 2 chips)	180,000	50	700	5	25	6.5	130,000	200 (250)	99.0–99.5	Homopolymers
ION G5-S5, P530 chip 600 SE	60,000	50	500	40	70	7	12,000	570 (650)	99.3–99.7	Homopolymers, low-quality ends
MGI Tech DNBSEQ-T7 (2 × 150 PE; maximum 4 cells)	990,000	50–100	6,000	1.2	4.5	24	5,000,000	250 (290)	99.9	Low-quality ends
MGI Tech DNBSEQ-G400RS (2 × 200 PE; 400 SE maximum 2 cells × 4 lanes)	480,000	50–100	1,800	4	11	108	450,000	350 (400)	99.9	Low-quality ends

a Costs are given in EUR. Information is compiled from multiple recent literature sources and price requests from sequencing companies. FC, flow cell; PE, paired-end; SE, single-end.

Tedersoo L et al Perspectives and Benefits of High-Throughput Long-Read Sequencing in Microbial Ecology. Appl Environ Microbiol. 2021;87: e0062621. doi:10.1128/AEM.00626-21

Short vs. Long Read Sequencing

- Short Read:
 - Cheaper (per base)
 - Higher capacity / throughput
 - Higher accuracy
 - Reads are consensus of many molecules
- Long Read:
 - More expensive (per base)
 - Lower capacity /throughput
 - Lower accuracy
 - Capable of single molecule sequencing

Pathogen Genomes

- Bacteria:
 - Typically contained within a **single** circular chromosome (some are linear)
 - **Haploid** genomes
 - May contain **plasmids** (extrachromosomal DNA)
 - Genome size range from 0.5Mb to ~10Mb (average is about 3-5Mb and contain about 3000-5000 genes)
- Viruses:
 - Can be DNA or RNA, single stranded or double stranded (classified into 7 families)
 - Range from 1-2Kb to ~1-2Mb (Pandoravirus salinus = 2.5Mb!)
 - Depend on host cellular mechanisms to replicate
- Eukaryotic Parasites (fungi, protists, and worms):
 - Usually a few to a few hundred Mbs
 - Usually multiple chromosomes

Microbial Genomes are Constantly Evolving – driving forces

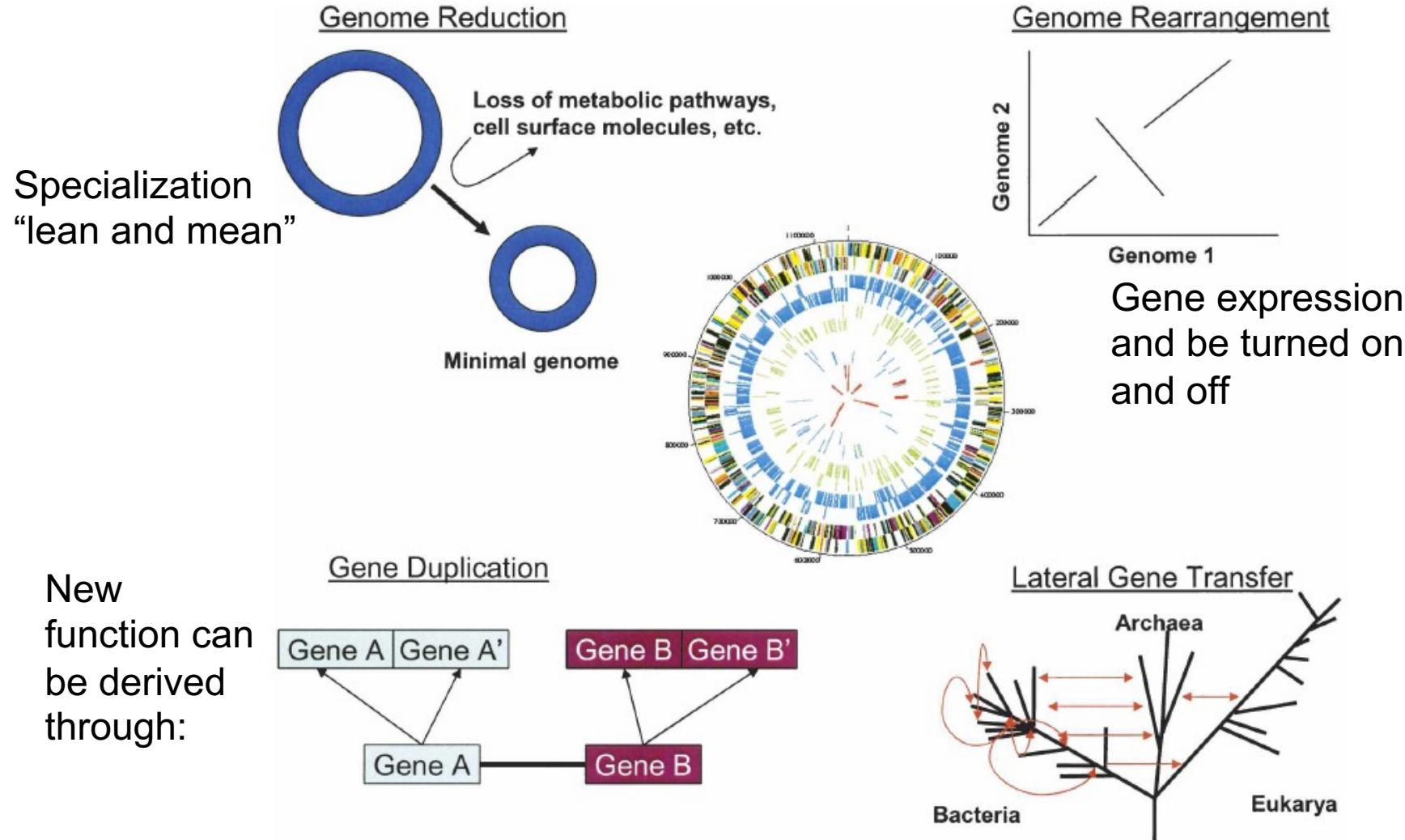
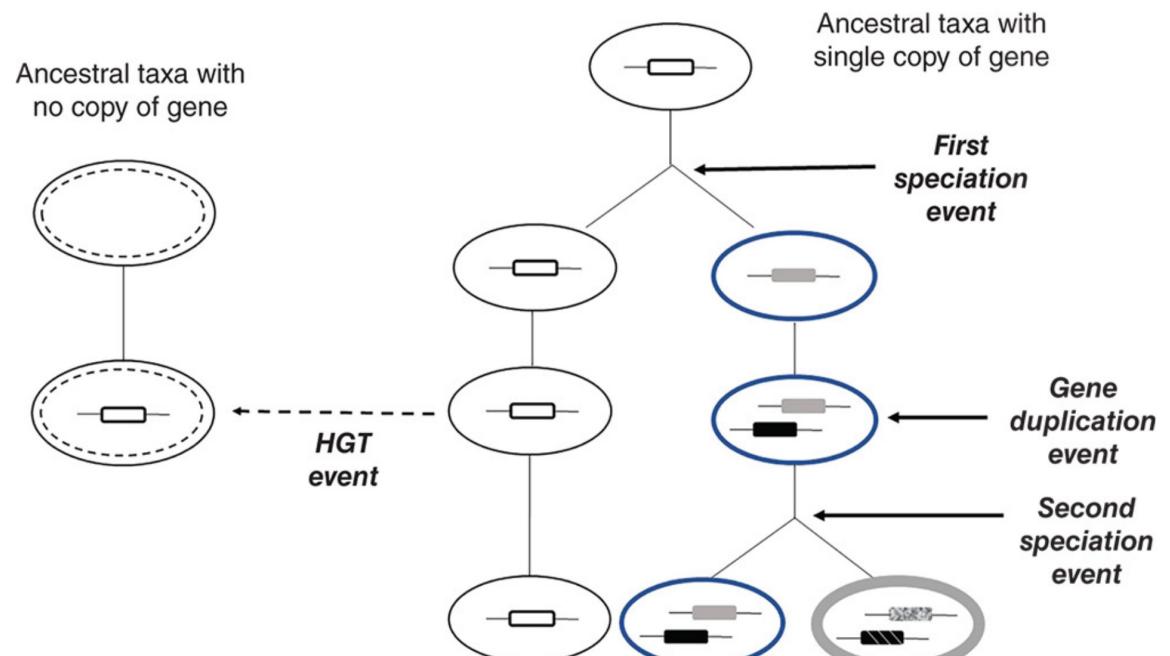


Figure 2. Multiple forces, including genome reduction, genome rearrangement, gene duplication, and acquisition of new genes via lateral gene transfer, are shaping microbial genomes. Details of each of these processes can be found in the text.

Homology

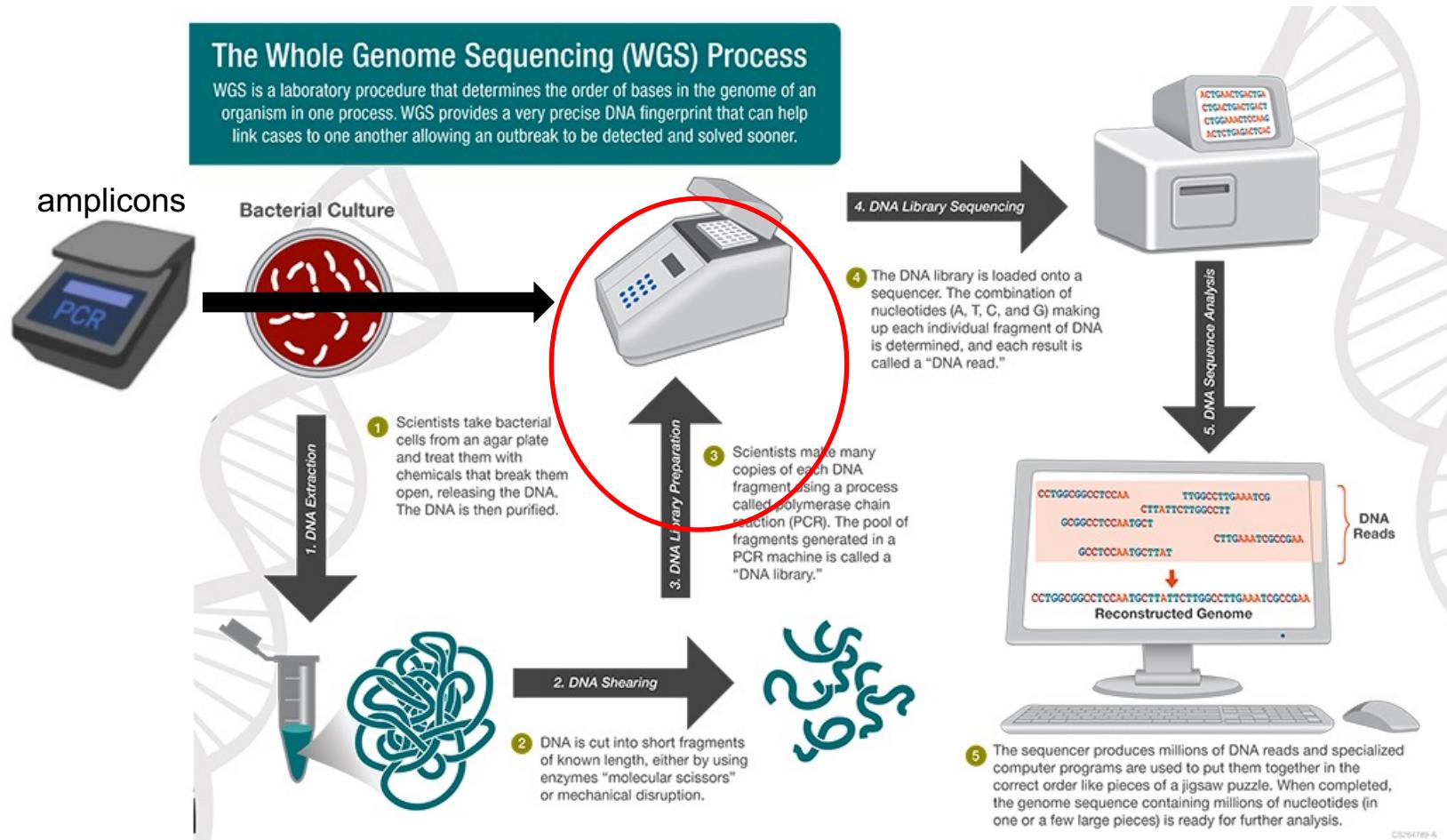
- Homology: similarity due to **shared common ancestry**
- **Orthologs**: arise due to speciation
- **Paralogs**: arise due to gene duplication
- **Xenologs**: arise due to horizontal gene transfer (no homology!)



If the black gene is deleted in one sister taxon and the gray gene is deleted in the other; this can result in mis-interpretation of the true relationship of the two genes

HGT event could also lead to misinterpretation of phylogeny

Whole Genome Shotgun Sequencing with NGS



<https://www.cdc.gov/pulsenet/pathogens/protocol-images.html#wgs>

Sequence Data Analysis

- Millions to Billions of partially overlapping reads were generated from a single “sequencing run”
- Steps of Data Analysis
 - Assembly
 - De-novo Assembly
 - Mapping
 - Annotation (adding information to sequences)
 - Gene Prediction (determine coding vs. non-coding regions of the genome)
 - Functional Prediction (determine the functions of genetic elements)
 - Variant Identification (single nucleotide variants, Allelic differences) – module 3

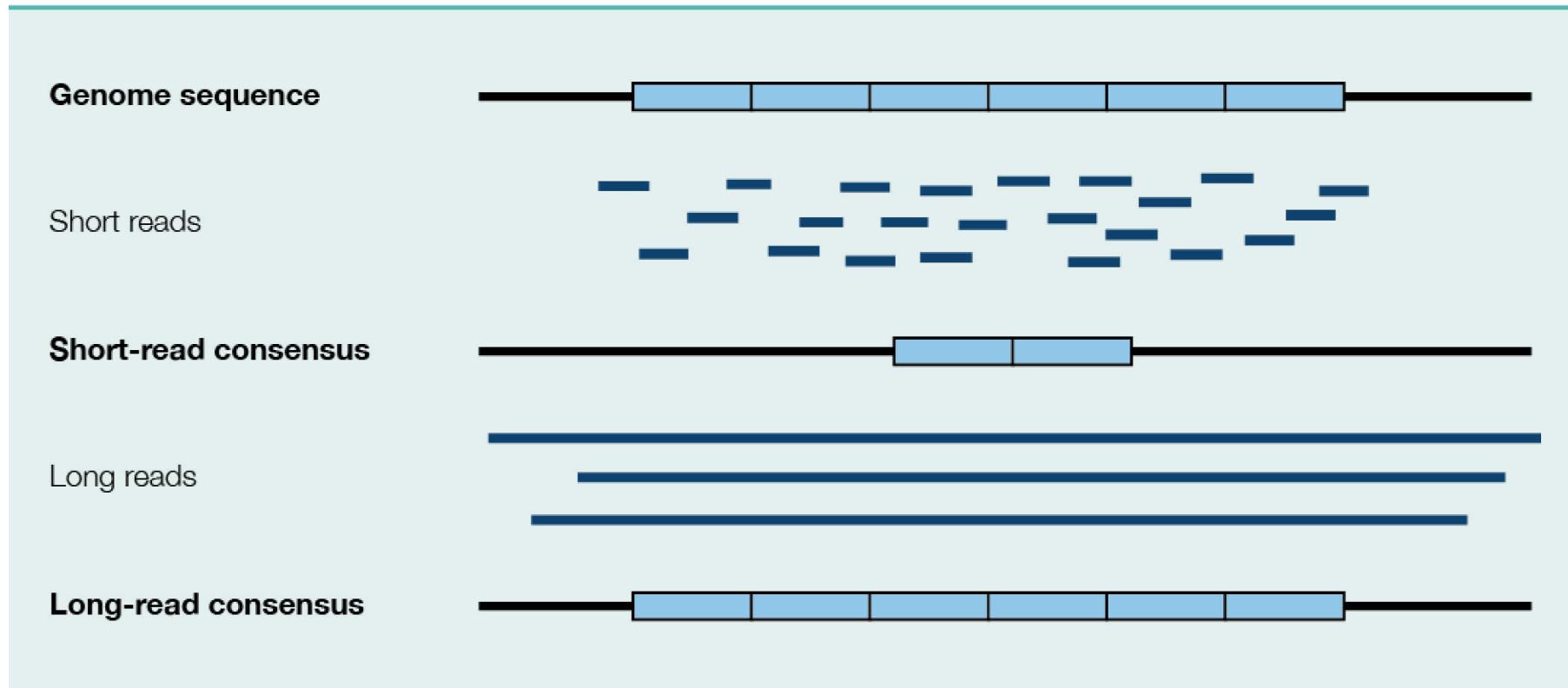
Genome Assembly

- The task is to reconstitute the whole genome from fragments of DNA sequences
- Two basic approaches:
 - **Denovo** assembly
 - Different computer algorithms available to identify overlapping sequences and merge them -> assemble
 - Reference Assisted assembly (aka **mapping**)
 - Map sequences to an existing related genome sequences
 - Assembly quality affected by the reference genome used (better mapping if the reference genome is closely related)



Assembly lecture from CBW: <https://www.youtube.com/watch?v=sysnKQvqmnk>

Why are long reads beneficial?



- *de novo* assembly of short read sequences – repetitive regions

NGS Error Rates

- Sequencing instruments have different error rates and are prone to different error types.
 - **Sanger** – prone to **substitution** errors and 0.1-1% error rate
 - **Ion Torrent** – prone to **indel** and 1% error rate
 - **SOLiD** – prone to **A-T bias** and >0.06% error rate
 - **Illumina** – prone to **substitution** errors and >0.1% error rate
 - **PacBio and MinION (3rd Gen Sequencer)** >5% error rate
- To minimize errors, same regions of the genome typically sequenced multiple times (>30X). This is called sequencing **depth coverage**. The consensus is taken as the correct sequence.

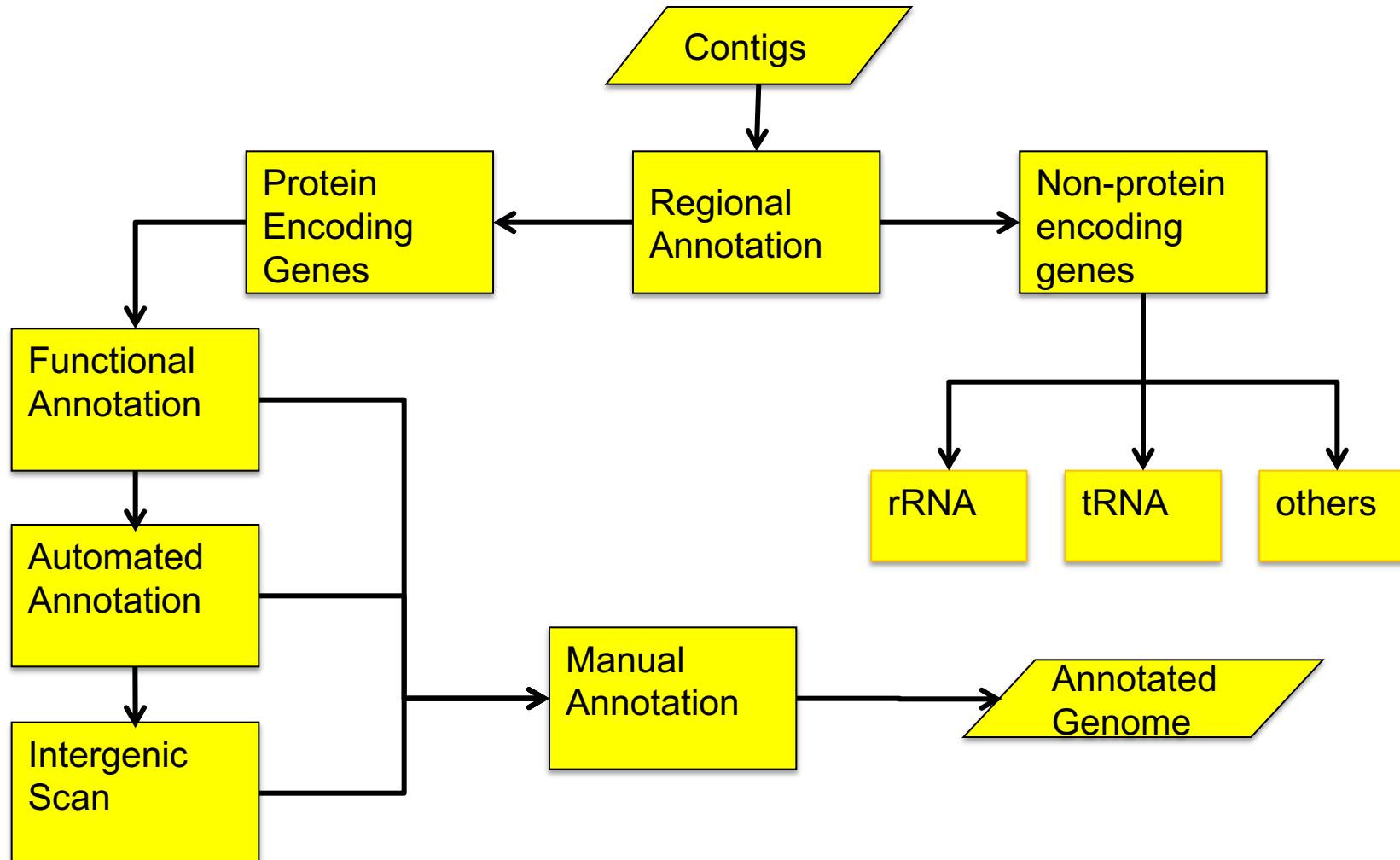
Contigs vs. Complete Genomes

- After assemblies, there are often still “gaps” in the genomes that can not be closed due to lack of sequence coverage or, more likely, un-resolved repeats
- So instead of the complete genome, you get a set of contiguous sequences (**contigs**) representing most of the genome
- Closing the gaps manually (aka “**finishing**” the genome) is labor intensive and expensive so very few groups still do that
- Long reads from 3rd generation sequencers can improve assembly

Genome Annotation

- Strings of A's, T's, C's and G's do not mean much on their own to us
- The goal of genome annotations is to identify **genes** and other features and locate them on the genomic sequence
 - Functions
 - Locations
- Locations of coding genes in a sequence can be identified by computer program because coding genes have different "word" frequencies compared to non-coding sequences – **ab initio gene prediction**
 - Works well for microbial genomes which are compact and usually no intron

Annotation Overview



Function Prediction

- The most common way to assign functions to a protein-coding gene is by **sequence similarity search**
- **It is assumed that genes that have sequence similarity are derived from the same ancestral gene and, therefore, have similar functions**
- BLAST is the most common tool for performing similarity search.
- Infer the function of one gene/protein based on its similarity to another gene/protein of known function is called "**transitive annotation**"
- Requires a database of known genes (e.g. GenBank)

BLAST

- Basic Local Alignment Search Tool
- Developed in 1990 and 1997 (S. Altschul)
- A heuristic method for performing local alignments through searches of high scoring segment pairs (HSP's)
- 1st to use statistics to predict significance of initial matches - saves on false leads
- Offers both sensitivity and speed
- Most highly cited bioinformatics tool!

Automated Annotation Systems

- Similarity searches form the foundation of annotation systems – thousands if not millions of searches are performed to annotation one genome
- Prokka (<http://www.vicbioinformatics.com/software.prokka.shtml>)
- Bakta (<https://github.com/oschwengers/bakta>)
- NCBI Prokaryotic Genome Annotation Pipeline
(http://www.ncbi.nlm.nih.gov/genome/annotation_prok/)

Some have a web submission form so you can submit your genome to be annotated for free!

Comparative Genomics

- The goal for comparative genomics is to identify **genomic variations** that can be correlated to phenotypic characteristics of an organism
- For example: we might be interested to know why certain isolates of a pathogen are more resistant to antibiotics or more virulent
 - E.g. comparing two strains of *E. coli* could reveal that one contains shiga toxins (bloody diarrhea) and the other does not (commensal)
 - We can also use these variations to track the transmission of pathogens

Comparative Genomics

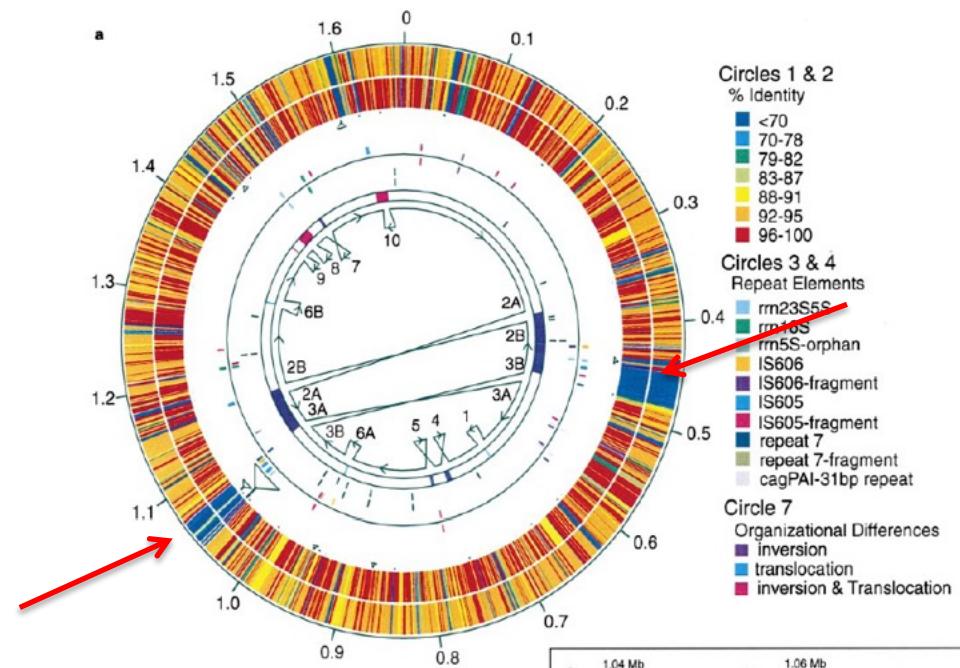
- Variations can occur at different levels
 - **Regional** (a stretch of DNA is present in one isolate but absent in another)
 - Strain-specific chromosomal regions
 - Strain-specific plasmids
 - VNTRs (tandem repeats)
 - **Gene** (a gene is missing or codes for different amino acids in one isolate compared to another)
 - Strain-specific genes
 - Allelic differences
 - Only looking at genes and not regulatory elements or non-coding genes
 - **Nucleic Acid** (a single nucleotide is different in one isolate compared to another)
 - Single nucleotide variants (SNVs) or single nucleotide polymorphism (SNPs)

First Comparative Genomics Paper

- published in 1999
 - 2 *Helicobacter pylori* genomes isolated 7 years apart were compared

Found more than half of the **strain specific genes** are clustered in hyper variable regions (red arrows)

This observation soon was consistently observed in many other species

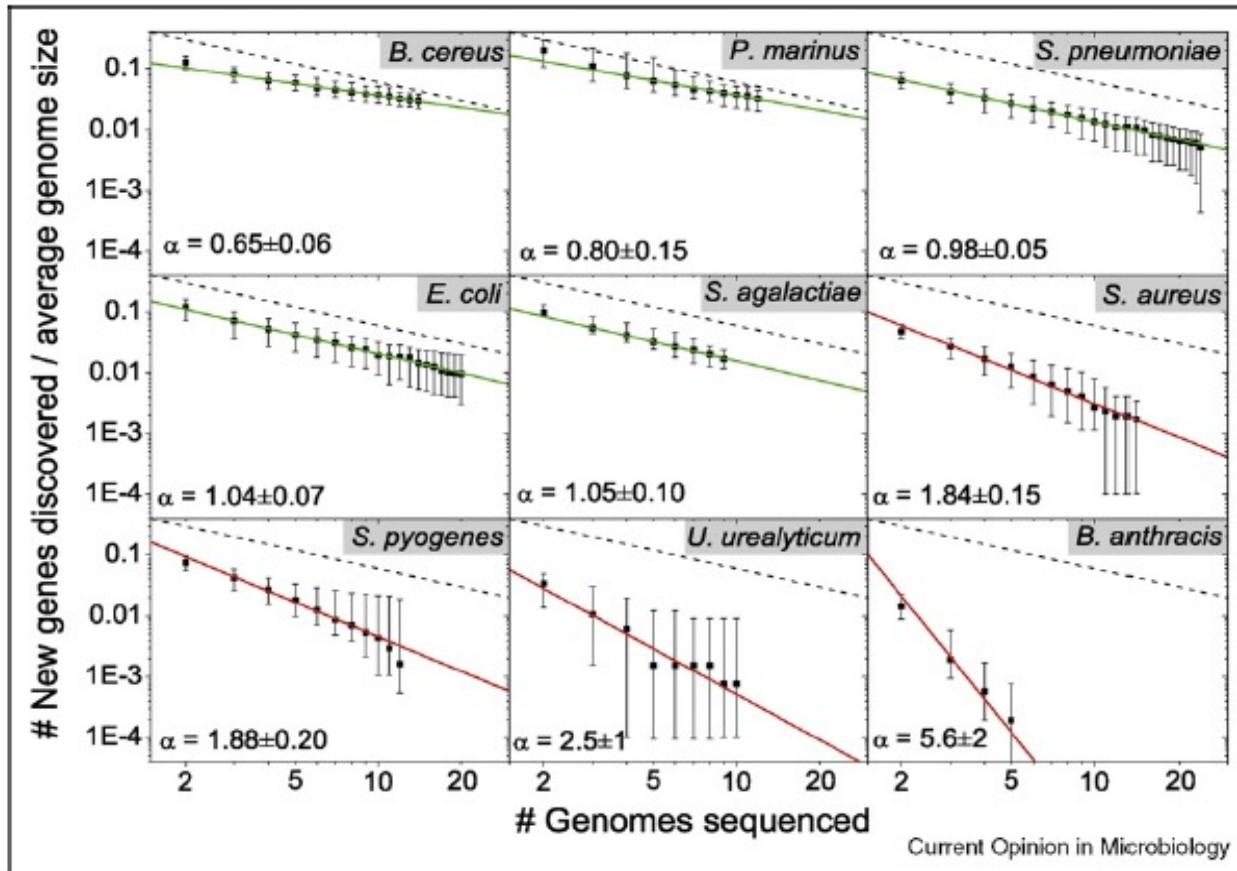


Alm et al, Nature 1999

Pan-genomes

- Comparative Genomics and huge genomic variations among strains in a bacterial species lead to the idea of pan-genomes
- The term first coined in 2005 in a paper by Tettelin et al., in which they compared sequenced genomes from six *S. agalactiae* (*group B Streptococcus*).
- Pan-genome is the entire gene set of a species - consists of the **core (housekeeping)** genes of strains in a species + **strain-specific (accessory)** genes
- Pan-genome calculation extrapolates observations based on a limited number of strains to come up with the theoretical number of genomes required to fully capture the pan-genome of a species

Open vs. Closed pan-genome



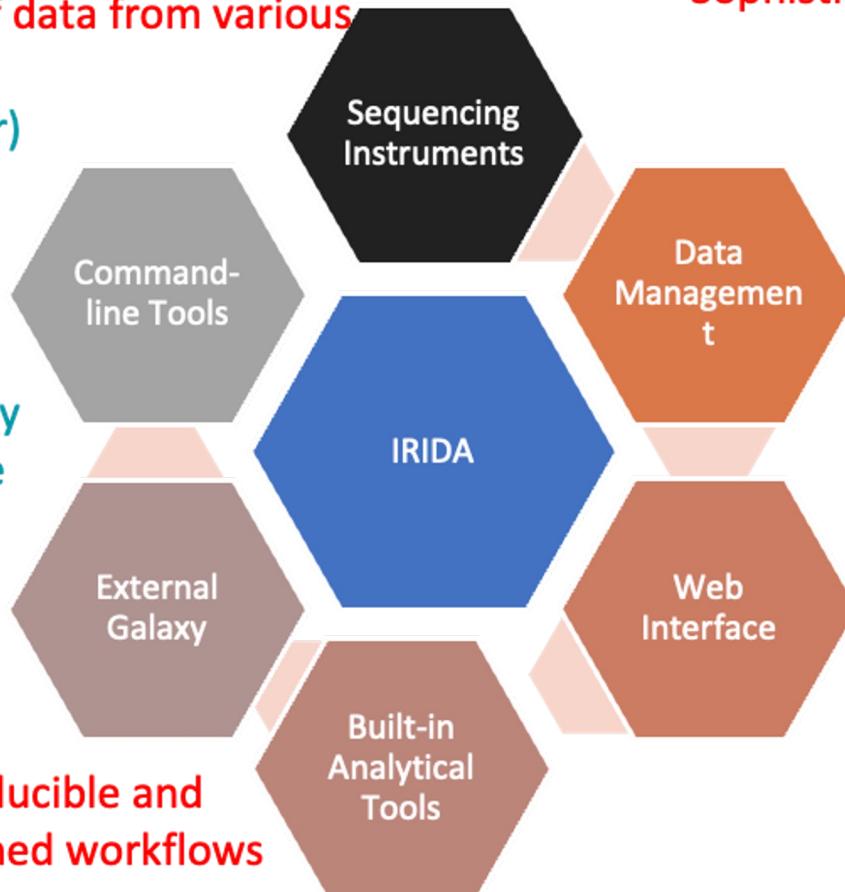
If the number of strains needed to capture the pangenome of a species is finite (red lines), then the species have a **closed pan-genome**; otherwise it has an **open pan-genome**.

This concept helps microbiologists predict how a species will evolve over time (e.g. is it likely to acquire many resistance genes?).

IRIDA: Canada's Genomic Epidemiology Analysis Platform

Easy, automated transfer of data from various sequencing platforms
(SeqUDAS & IRIDA Uploader)

Tiered Analytic Platforms:
Data can be export seamlessly to Galaxy or to command line for custom analyses



Expanding Analytical Workflows with plug-in architecture (Assembly, annotation, SNP, MLST, AMR, serotyping, and more.)

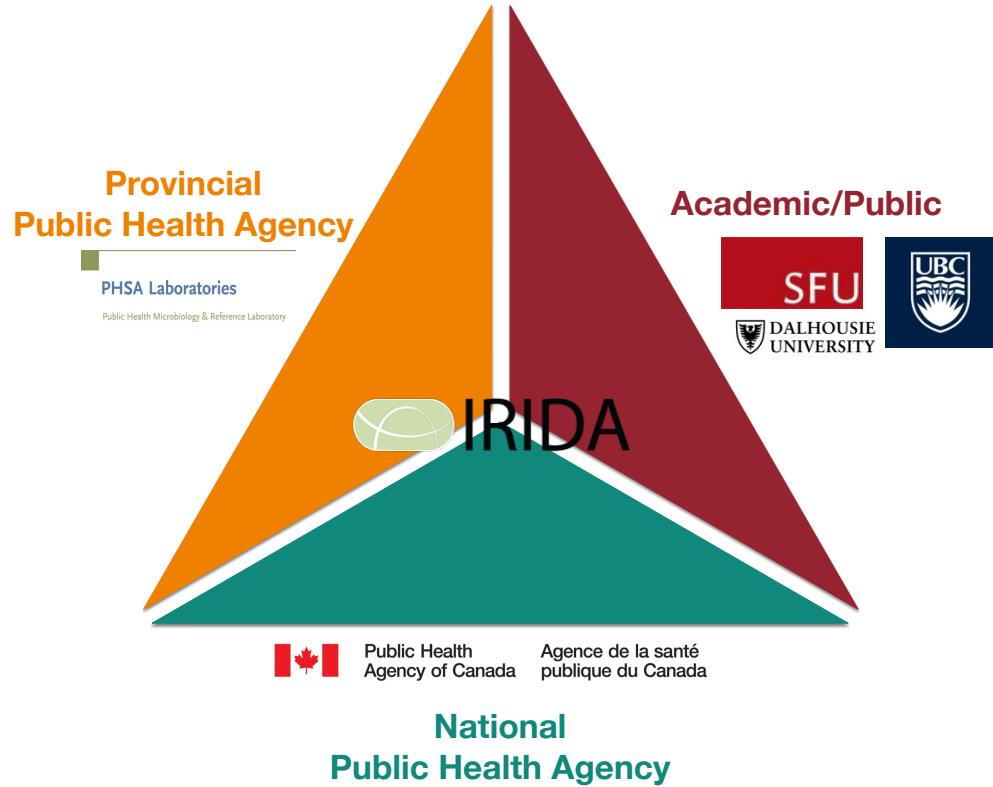
Sophisticated Project and Sample Management

ID	Name	Organization	Samples	Created	Modified
48	Mot_Salmonella_Feb2018	Salmonella enterica subsp. enterica	115	Feb 8, 2018 5:05 PM	May 2, 2018 2:31 PM
76	Mot_Salmonella	Salmonella	115	May 2, 2018 10:57 AM	May 2, 2018 1:12 AM
74					

Simple User Interface

Ready to Launch?

IRIDA is a Partnership

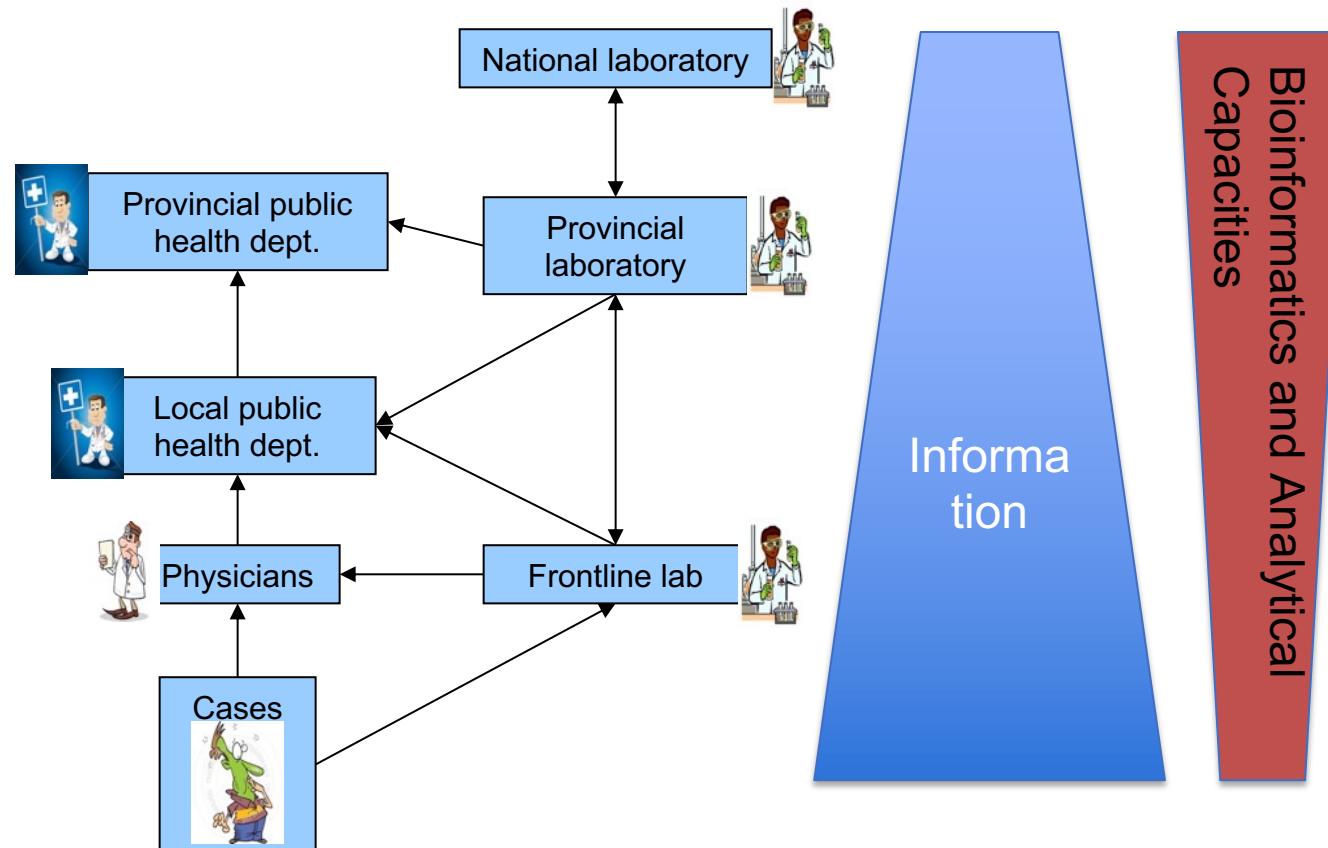


- Project Team has direct access to state of the art research in academia
- Project Team is directly embedded in user organization
- Interview with key stakeholders to understand the gaps and challenges

www.irida.ca

bioinformatics.ca

Challenges in Data Integration and Sharing - there are many players/systems



Modern software development



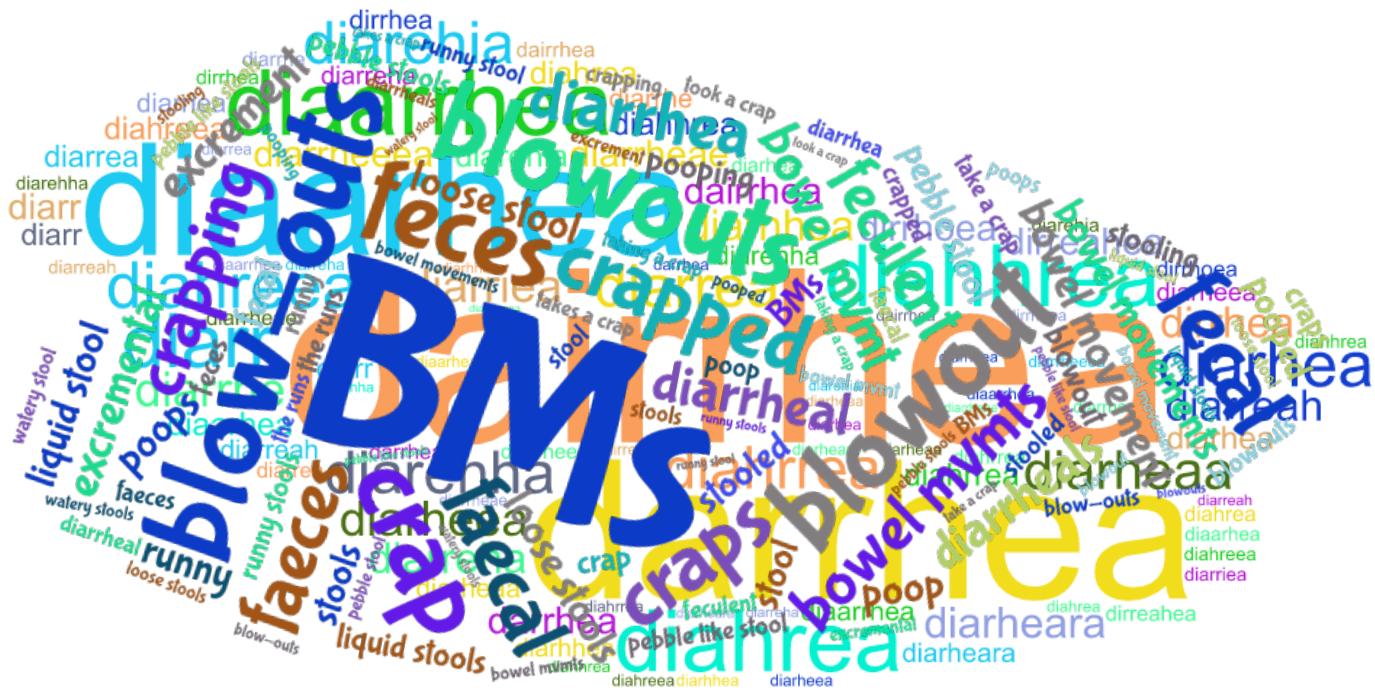
Contextual Information – Often Institutional Specific

- Examples:
 - Different acronyms or codes used for the same antibiotics
 - Different terminologies to describe the same concept or concepts with subtle differences (alcoholism vs. substantial alcohol use/abuse; died vs death, etc.)
 - Different units of measure
 - Different severity gradients



Metadata Problem: Spellings, Synonyms, Semantics

Diarrhea: “Abnormally increased frequency of loose or watery bowel movements.”



Fecal: "Portion of semisolid bodily waste discharged through the anus."

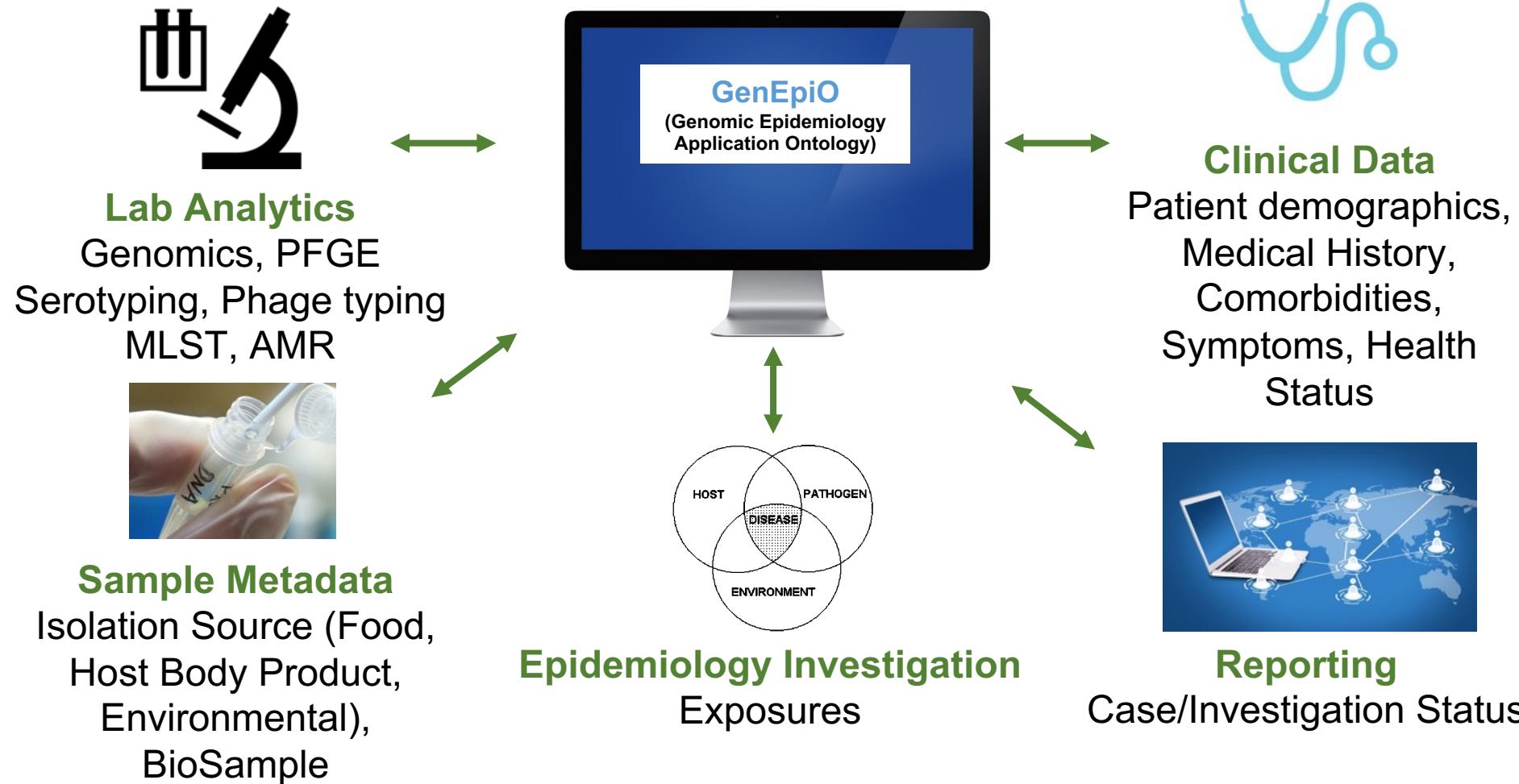
Misspellings and synonyms supplied by <http://project-emerse.org/>

Ontology

- A mechanism to specify and express a body of knowledge
- Standardized, well-defined **hierarchy of terms**
- Each term has a unique universal ID
- Terms interconnected with **logical relationships**
- Have formats that are Human AND computer readable
- This internally coherent tool can act as a universal translator of different standards



Gen-Epi-O: Combining Epi, Lab, Genomics and Clinical Data Fields



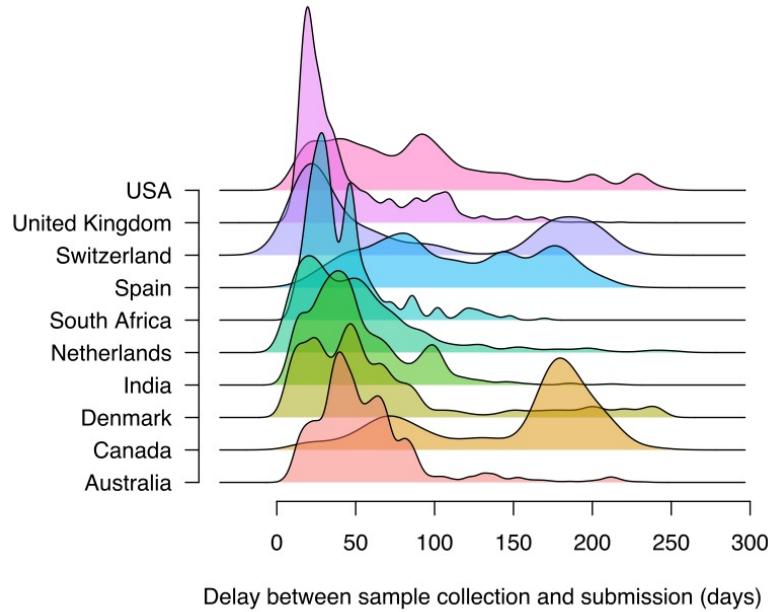
Challenges with Data Sharing in Canada

- Canada is comprised of **14 distinct healthcare** systems
- **No universal standard** for data collection or sharing (yet)
- No legally binding public health **data sharing agreement** in Canada
- No standard for COVID-19 patient data sharing between the **six BC health authorities**
- Arguably, the restrictive flow of patient data across Canada is a violation of the **universality** and **portability** of healthcare

Attaran, A. & Houston, A. Pandemic Data Sharing: How the Canadian Constitution Turned Into a Suicide Pact. (2020). doi:10.2139/ssrn.3612825

Delayed Data Release from Canada

To date (Mar2021), roughly 57,000 genomes sequenced and ~50% uploaded to GISAID (global repository of influenza and COVID consensus sequences)



Art Poon (Western): https://twitter.com/art_poon

Canada's missing SARS-CoV-2 genomes

Each circle's area is scaled in proportion to the number it represents.

On December 22, 2020, the Globe and Mail reports the Public Health Agency of Canada is examining over 25,000 SARS-CoV-2 sequences for evidence of lineage B.1.1.7.

By that date, 3,605 genomes from Canada have been released to GISAID, a global initiative to provide rapid and open access to pandemic data.

2,515 have full dates of sample collection

Only 337 genomes are released within two months of sample collection.

By December 22nd, about 520,000 cases have been confirmed in Canada.

visualized in R by @art_poon

Reasons cited for the delays in Canada

- Capacity to process the sequence data and metadata for public release from public health laboratories is low
 - E.g. dehosting of human sequences, de-identify and process the metadata
- Multiple sign-offs are needed for the release of data
- Privacy concerns barring the release of specific data with unique combination of metadata fields
- Wanting to ensure sequence data quality
- All these challenges are replicated across multiple jurisdictions - solving the same problems and implementing the solutions many times



Addressing Privacy Concerns in Sharing Viral Sequences and Minimum Contextual Data in a Public Repository During the COVID-19 Pandemic

Lingqiao Song^{1†}, Hanshi Liu^{1*†}, Fiona S. L Brinkman², Erin Gill², Emma J. Griffiths³, William W. L Hsiao², Sarah Savić-Kallesøe², Sandrine Moreira⁴, Gary Van Domselaar⁵, Ma'n H. Zawati¹ and Yann Joly¹

Open access

Original research

BMJ Open Canadians' opinions towards COVID-19 data-sharing: a national cross-sectional survey

Sarah A Savic Kallesoe ,^{1,2} Tian Rabbani ,^{1,3} Erin E Gill,⁴ Fiona Brinkman,⁴ Emma J Griffiths,¹ Ma'n Zawati,⁵ Hanshi Liu,⁵ Nicole Palmour,⁵ Yann Joly ,⁵ William W L Hsiao¹

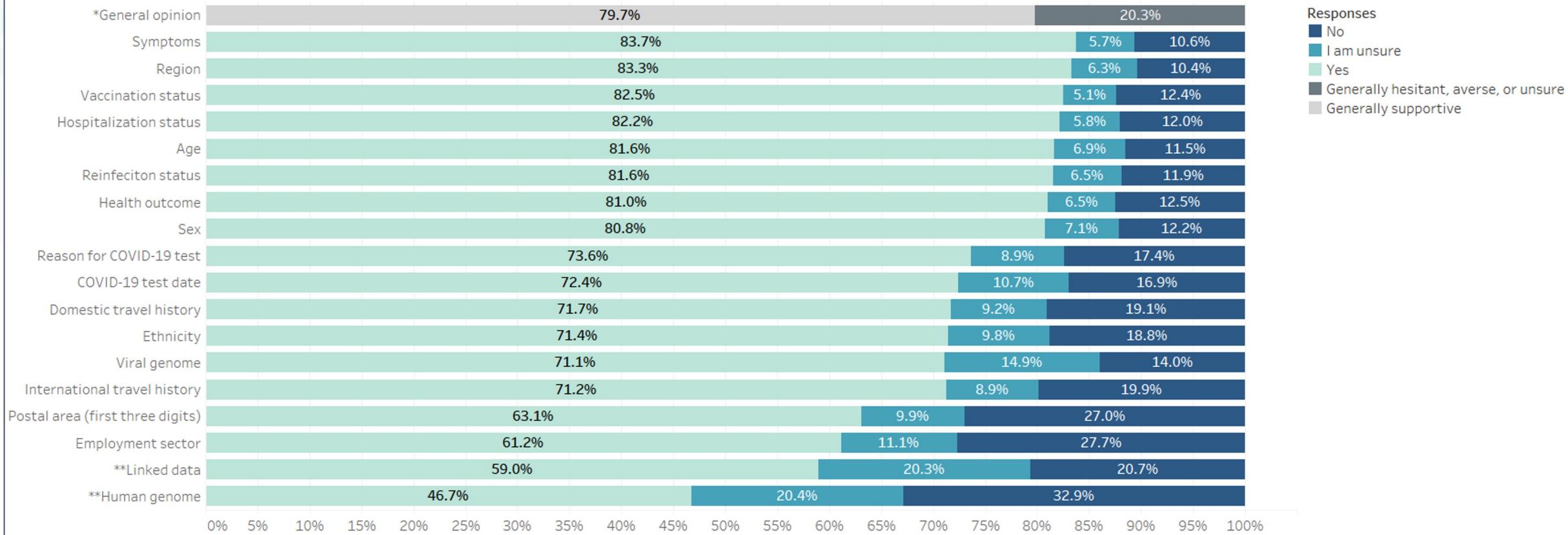
“What are Canadians’ opinions of public health authorities sharing anonymized COVID-19 data for research?”

Methods

- Cross-sectional survey of a sample representative of the Canadian population (by age, sex, ethnicity, and region)
 - n = 5,014 (n = 4,981 responses with complete demographic data [age, sex, province])
- Sampling method: as close to simple-random sampling as financially feasible.
 - Partnered with a Canadian market research team, Leger, to recruit participants
- Community partners and pilot survey
 - Ran a pilot survey in Summer 2021, 10+ hours of community partner feedback, cross-institutional partnerships

Canadians favour de-identified public health data sharing

Would you be comfortable with the following anonymized COVID-19 data collected from the population by public health authorities, which could potentially include your data, being publicly accessible?



*Participants' responses to the 16 datatypes are summarized in the "General opinion" variable. Participants who responded "Yes, I would be comfortable with this anonymized datatype being publicly shared" to nine or more datatypes were classified as "generally supportive". Those who responded "Yes" to eight or less were classified as "generally hesitant, averse, or unsure". Participants' responses to their comfort on human genome and linked data being shared with authorized researchers are not included in this variable.

**Only accessible to authorized researchers and not the public. Authorized researchers are those who have been approved to use the data and agreed, in writing, not to attempt to uncover the identity of the person or to share the data with unauthorized third parties.

Pyramid of Context

National surveillance

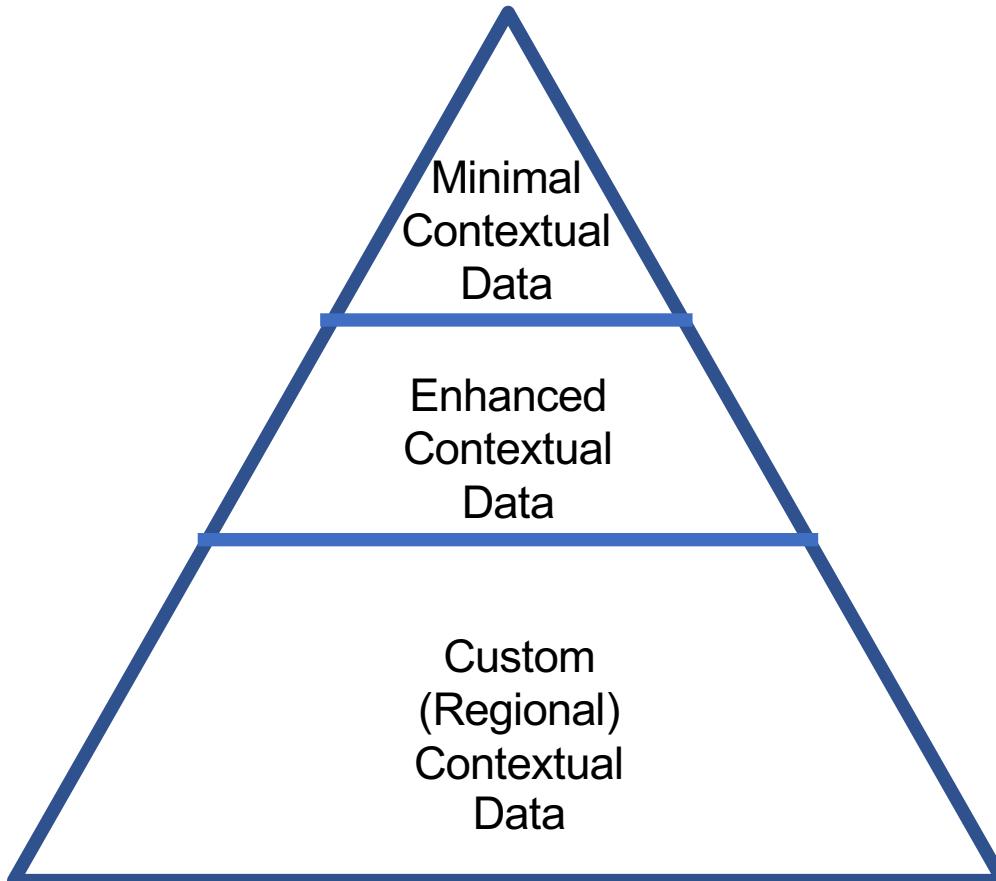
- Collection date
- Province
- Age
- Gender

Beyond minimal set

- Research/analyses

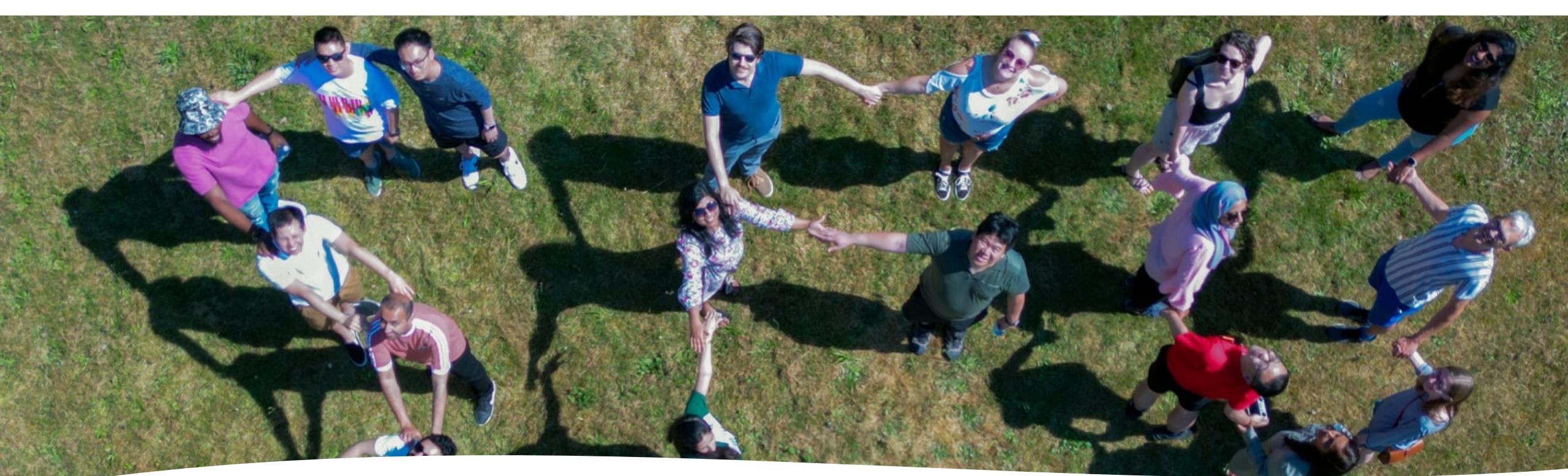
Province-specific

- Could be shared or stored locally



Minimal Contextual Data

Specimen collector sample ID
Sample collected by
Sequence submitted by
Sample collection date
Geo_loc name (country)
Geo_loc name (province/territory)
Organism
Isolate
*Isolation source
Host (scientific name)
Host disease
Host age
Host gender
Sequencing instrument
Consensus sequence software name
Consensus sequence software version



Thank you for your attention

WWW.CIDGOH.CA

bioinformatics.ca

We are on a Coffee Break & Networking Session

Workshop Sponsors:



Canadian Centre for
Computational
Genomics

