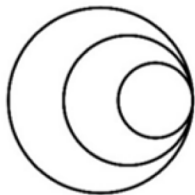


How to design and deliver pathogen genomics training for health and research professionals

Module 3C
Data Analysis and Integration
07/03/22
Silvia Argimón



**wellcome
connecting
science**



Centre for Genomic
Pathogen Surveillance



Session Outline

Aim: How to communicate to trainees the importance of the integration of different sources of data for decision making

- (Genomic) surveillance objectives
- Data sources and challenges
- Data analysis and integration
- CGPS tools for data integration and interactive visualization
- Activity and discussion

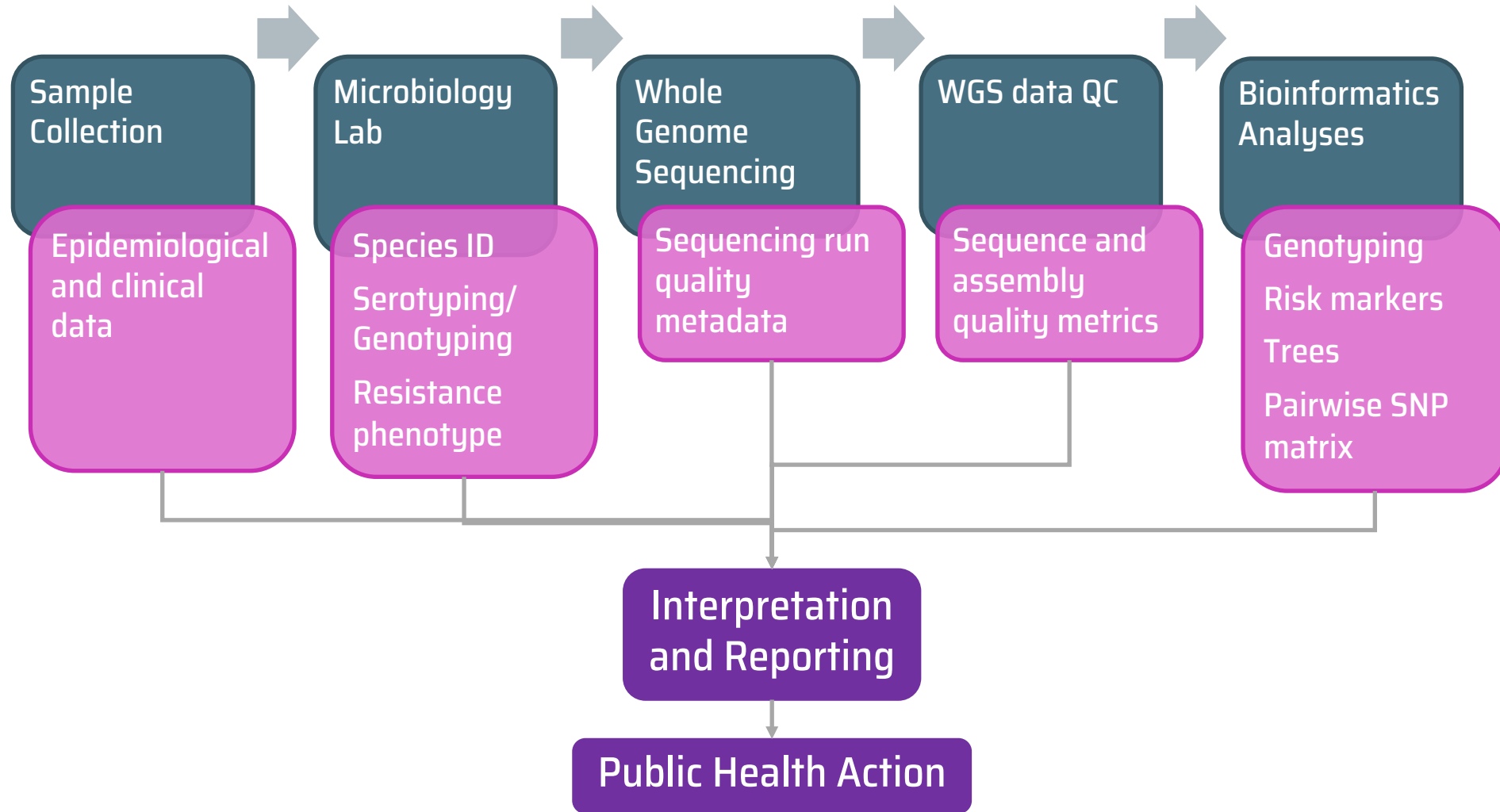
Surveillance Objectives for Infection Control

- Identifying and mapping of **high-risk clones** (HiRiCs) and **high-risk elements**
- Understanding the **routes of transmission** and the role and behaviour of **vectors** in the dissemination of HiRiCs
- Understanding the role of different **sources** or **reservoirs** in the dissemination across environmental, animal, and human habitats

Hajo Grundmann (2014). Towards a global antibiotic resistance surveillance system: a primer for a roadmap, *Upsala Journal of Medical Sciences*, 119:2, 87-95



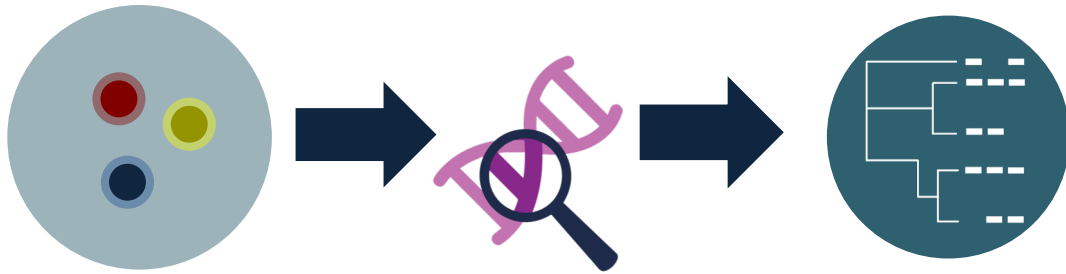
Genomic Surveillance - Data Sources and Integration



Data Challenges 🤯

- **Completeness and consistency**
 - Collecting hospital doesn't always record the patient's diagnosis
 - Different labs might test different panels of antibiotics
- **Interoperability**
 - Different computer platforms
- **Standardization**
 - Vietnam vs Viet Nam, or 07/10/2019 vs 10/07/2019
 - CLSI, EuCAST, BSAC susceptibility breakpoints
- **Different end-users**
 - Physicians, infection control, hospital authorities, Ministry of Health
- **Formats**
 - Output of one tool is not in the required format to use as input of the next tool
 - Integrating data from different sources (.pdf, .xlsx, .csv, .tre)

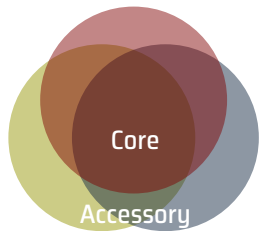
Genomic Surveillance - Data analysis



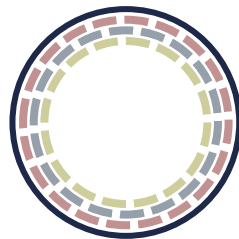
Extracting epidemiologically relevant information from the genomes

- Genetic relatedness
 - Genotyping (MLST, pathogen-specific)
 - Phylogenetic relationships (tree)
- Presence of markers of risk

Pangenome Analysis



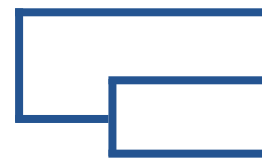
Mapping to a Reference



1 AAT**C**GCTTTACGACCAG...
 2 AATGGCTTTAT**T**GACAAG...
 3 AAT**C**GCTTTAT**T**GACAAG...
 * * *

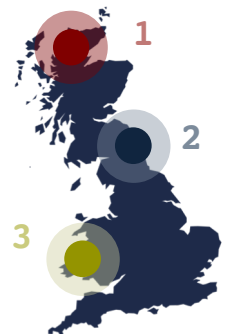
BLAST assemblies  map reads

SNP-based tree

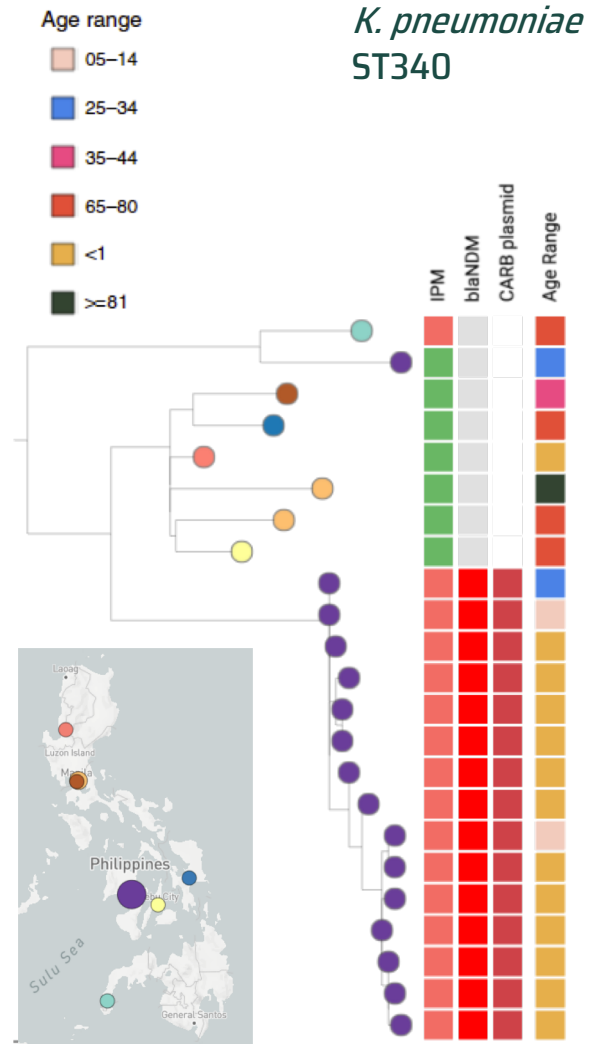


Newick format: ((2,3)1);

	gene1	gene2	gene3
1	■	□	□
2	□	■	■
3	■	■	□

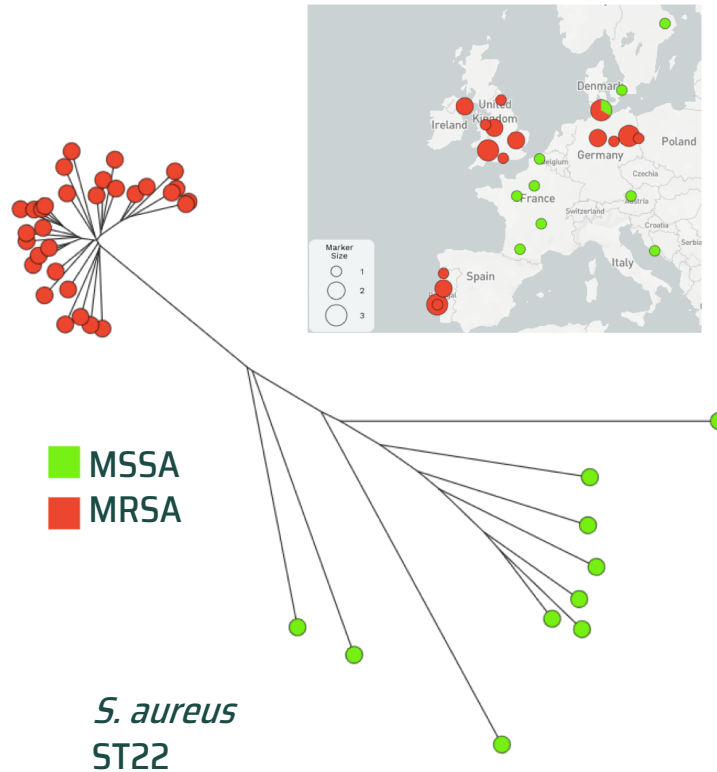


High-risk clones and genomic epidemiology



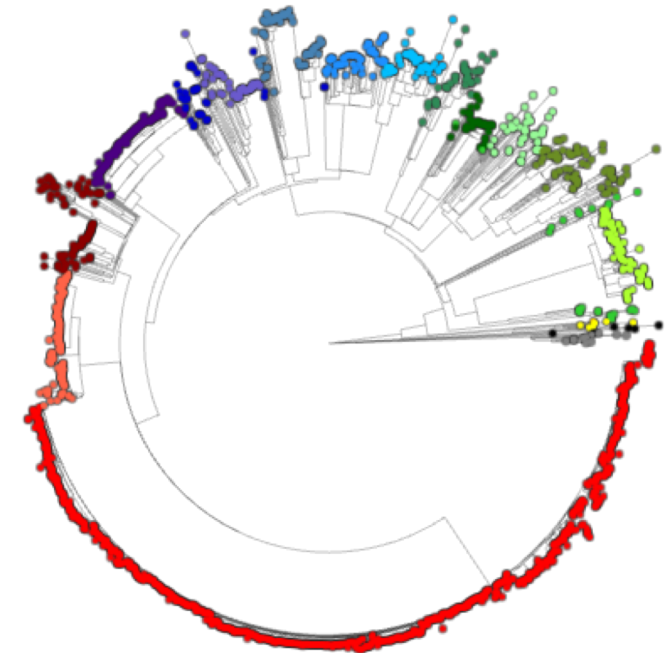
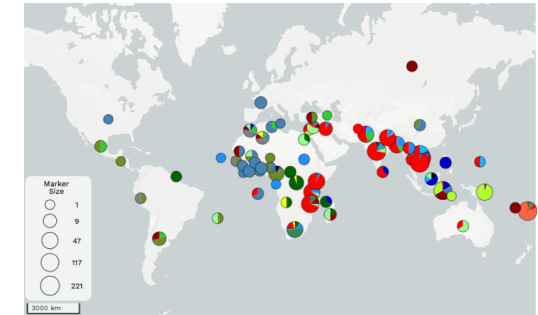
Argimón S. et al. 2020. Nat Comms 11:2719

- Clonal Relatedness
- Abundance
- Spatial/temporal distribution
- Risk properties



Aanensen D. et al. 2016. Mbio 7(3):e00444-16

S. Typhi
4.3.1 (H58)



Wong V. et al. 2015. Nat Genet 47:632

CGPS - Free Web Applications for Pathogen Surveillance

Collection



epi**collect**

<https://five.epicollect.net>

Integration

Data-flo

<https://data-flo.io/>

Interactive Visualization



Micro**react**

<http://microreact.org>

Analysis



Pathogen**watch**

<https://pathogen.watch>

@Pathogenwatch @MyMicroreact @EpiCollect



wellcome
connecting
science



T3**connect**



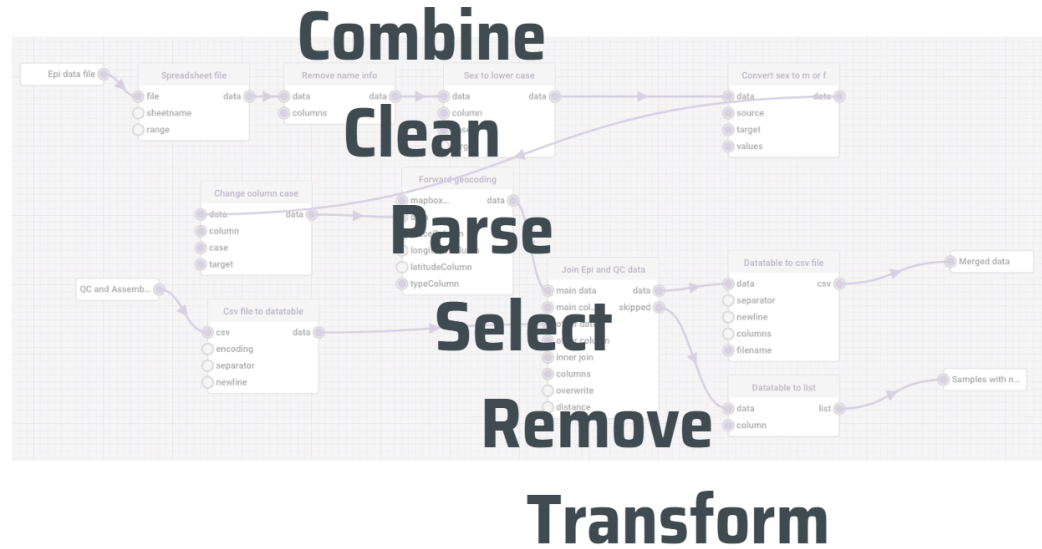
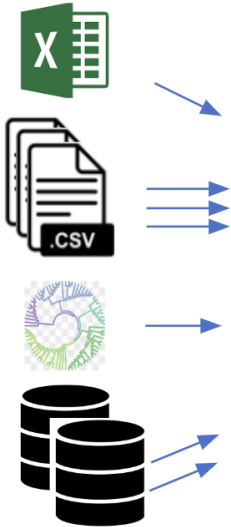
Centre for Genomic
Pathogen Surveillance



<https://data-flo.io/>

Customised integration and manipulation of diverse data via a simple drag and drop interface

SOURCES



DESTINATIONS

- Microreact
- Google Drive
- CSV
- TSV
- SQLite DB file
- Etc etc



Interactive visualization of clustering (trees), geographic (map) and temporal (timeline) data.

Whole genome sequencing reveals potential spread of *Clostridium difficile* between humans and farm animals in the Netherlands, 2002 to 2011

C W Knetsch¹, T R Connor², A Mutreja³, S M van Dorp¹, I M Sanders¹, H P Browne³, D Harris¹, L Lipman⁴, E C Keessen⁴, J Corver (j.corver@lumc.nl)⁵, E J Kullper¹, T D Lawley³

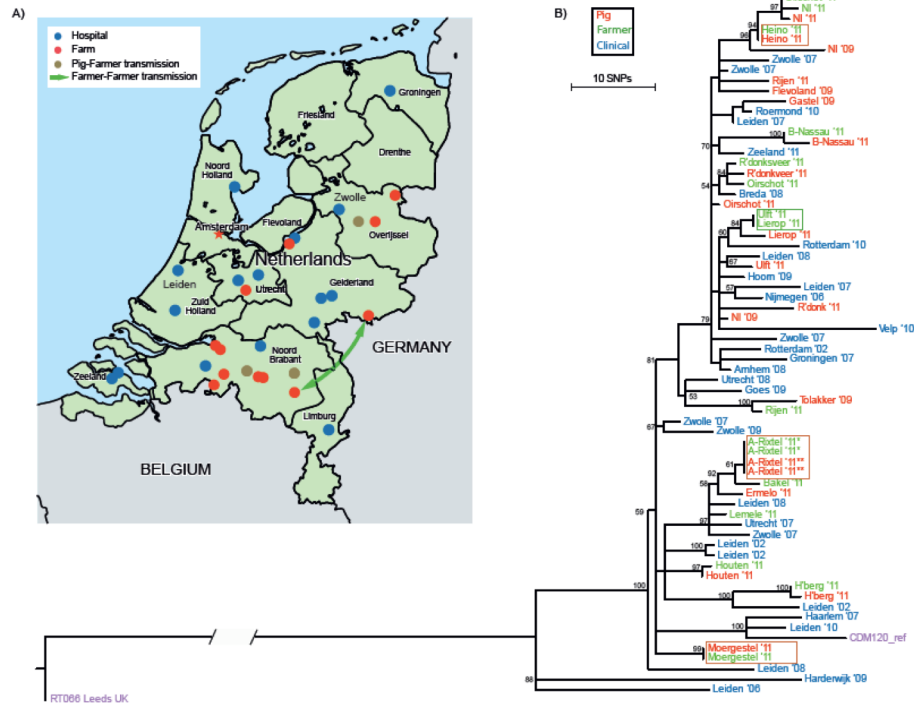


TABLE 1A

Clostridium difficile type 078 isolates used in this study, the Netherlands, 2002–11 (n=65)

R_L#T ^a	Year	City	RT	Isolate	Source	Related isolates	Association	ENA ID ^b
8080_2#24	2006	Leiden	078	6072310	Clinic	Non-outbreak	Healthcare	ERS138026
8080_2#25	2006	Nijmegen	078	6086336	Clinic	Non-outbreak	Healthcare	ERS138027
8080_2#26	2007	Leiden	078	7001233	Clinic	Non-outbreak	Healthcare	ERS138028
8080_2#27	2007	Groningen	078	7004578	Clinic	Non-outbreak	Unknown	ERS138029
8080_2#28	2007	Utrecht	078	7005405	Clinic	Non-outbreak	Unknown	ERS138030
8080_2#29	2007	Zwolle	078	7021455	Clinic	Non-outbreak	Healthcare	ERS138031
8080_2#30	2007	Zwolle	078	7044912	Clinic	Non-outbreak	Community	ERS138032
8080_2#31	2007	Zwolle	078	7066827	Clinic	Non-outbreak	Community	ERS138033
8080_2#32	2007	Zwolle	078	7071308	Clinic	Non-outbreak	Healthcare	ERS138034

FIGURE 2

Phylogenetic cluster showing relatedness of *Clostridium difficile* clinical, pig and farmer isolates, the Netherlands, 2008–11 (n=4)

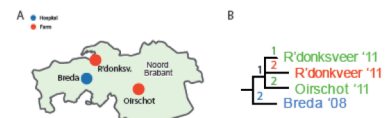


FIGURE 3

Phylogeny of *Clostridium difficile* 078 isolates showing the presence of antimicrobial resistance determinants, the Netherlands, 2002–11 (n=65)

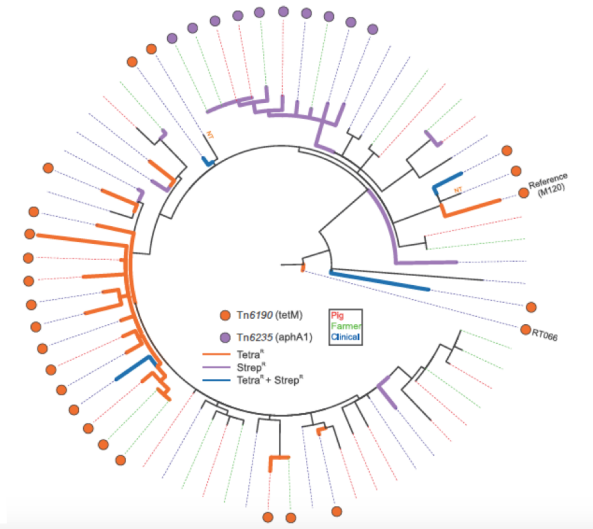


TABLE 3
Results of Antimicrobial susceptibility testing

Isolate	Source	Strepto- mycin	Tn6235	Tetra- cycline	Tn6190
6072310	Clinic		Absent		Present
6086336	Clinic		Absent		Present
7001233	Clinic	NT	Absent		Absent
7004578	Clinic		Absent		Present
7005405	Clinic		Present		Absent
7021455	Clinic		Present		Absent
7044912	Clinic		Absent		Present
7066827	Clinic		Absent		Absent
7071308	Clinic		Absent		Absent
7086074	Clinic		Absent		Absent
7091912	Clinic		Absent		Absent
801064	Clinic		Absent		Absent
8013820	Clinic		Absent		Absent
8061728	Clinic		Absent		Present
8053444	Clinic		Present		Absent
1129280	Clinic		Absent		Absent
H205	Farmer		Present		Absent
B37.3	Pig		Present		Absent
53737	Clinic	NT	Present	NT	Present
473737	Clinic		Absent		Present
H102	Farmer		Absent		Absent
B31.3	Pig		Absent		Absent
B37.3	Pig		Absent		Absent
H121	Farmer		Absent		Present
B27.7	Pig		Absent		Absent
7091912	Farmer		Absent		Present
H189	Farmer		Absent		Absent
B23.6	Pig		Absent		Present
H206	Farmer		Present		Absent
B15.1	Pig		Absent		Absent

Basic concepts in data analysis and integration

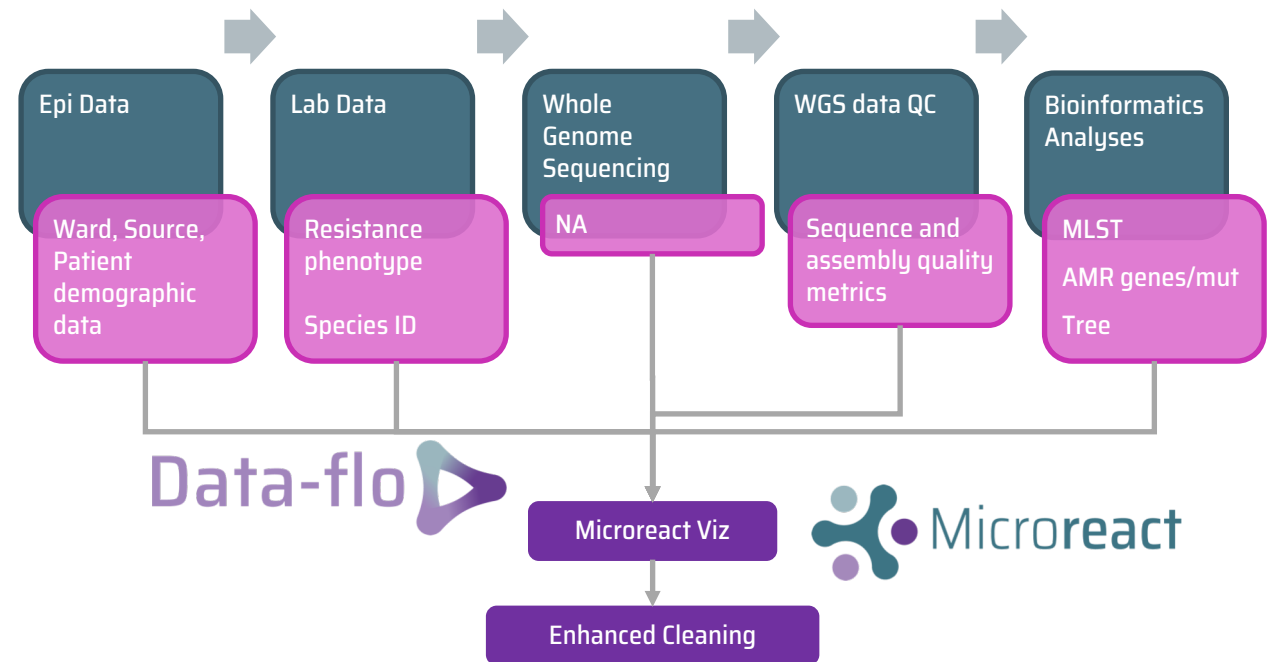
Domain Concepts	Strategies	Assessments	Resources
Specific analytic tool/pipeline	Demo/tutorial with dummy data. Presentation with link to tools/resources	Exercise with “real data”	Module 1B https://lms.welcomeconnectingscience.org/mod/forum/view.php?id=2525
Different types of data for genomic surveillance	Presentation Poll Group Activity	Wrap-up discussion	https://docs.data-flo.io/using-data-flo/data/data-types
Importance of data harmonization	Presentation (Bonus) Group Activity	Wrap-up discussion	https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8847733/
Data integration and visualization	Presentation Poll Group Activity	Wrap-up discussion	data-flo https://docs.data-flo.io/introduction/readme Microreact https://docs.microreact.org/ Tidyverse https://github.com/rstudio/cheatsheets/blob/main/data-transformation.pdf Nextstrain https://nextstrain.org/ Phandango https://jameshadfield.github.io/phandango/#/
Genomic surveillance/epidemiology	Presentation (Examples) Group Activity	Wrap-up discussion Interpretation of case study	https://alliblk.github.io/genepi-book/ https://www.futurelearn.com/courses/pathogen-genomics-a-new-era-in-global-health-surveillance-and-strategy https://www.futurelearn.com/courses/genomics-covid-19

Module 3C Exercise

Data analysis and integration for decision making

Hospital outbreak investigation → Enhanced cleaning

- Work in pairs
- Read the Background
- Activity 1. Data sources
- Activity 2. Integration and Viz
- Activity 3. Reflection (table groups)
- Session wrap-up.



Module 3C Wrap-up

Aim: Highlight the importance of the integration of different sources of surveillance data for decision making

Presentation (passive learning)

Activities (active learning)

Which domain concepts covered in this session (or similar ones) would you foresee having to teach in the future? How would you do it? Note: See presentation slides 2 and 11 for reference to domain concepts.

What do you think are the pros and cons of using interactive web tools such as data-flo and Microreact for teaching?

Can you think of alternative ways to teach this module?



References

Hajo Grundmann. 2014. Towards a global antibiotic resistance surveillance system: a primer for a roadmap. *Upsala Journal of Medical Sciences*, 119:2, 87-95
doi:[10.3109/03009734.2014.904458](https://doi.org/10.3109/03009734.2014.904458)

Argimón S, Abudahab K, Goater R, Fedosejev A, Bhai J, Glasner C, Feil E, Holden M, Yeats C, Grundmann H, Spratt B, Aanensen D. 2016. Microreact: visualizing and sharing data for genomic epidemiology and phylogeography. *M Gen*, 2(11): doi:[10.1099/mgen.0.000093](https://doi.org/10.1099/mgen.0.000093)

Knetsch CW, Connor TR, Mutreja A, van Dorp SM, Sanders IM, Browne HP, Harris D, Lipman L, Keessen EC, Corver J, Kuijper EJ, Lawley TD. 2014. Whole genome sequencing reveals potential spread of *Clostridium difficile* between humans and farm animals in the Netherlands, 2002 to 2011. *Euro Surveill*, 19(45):20954. doi: 10.2807/1560-7917.es2014.19.45.20954

Acknowledgements

This course was developed by a collaboration between the [Centre for Genomic Pathogen Surveillance](#) and [Wellcome Connecting Science](#). It was brought to you by [T3Connect – Data Science and Genomic Pathogen Surveillance Training Programme](#), funded by [UKRI](#).

This module contains materials from the following sources:

- [Storyset | Customize, animate and download illustration for free](#)

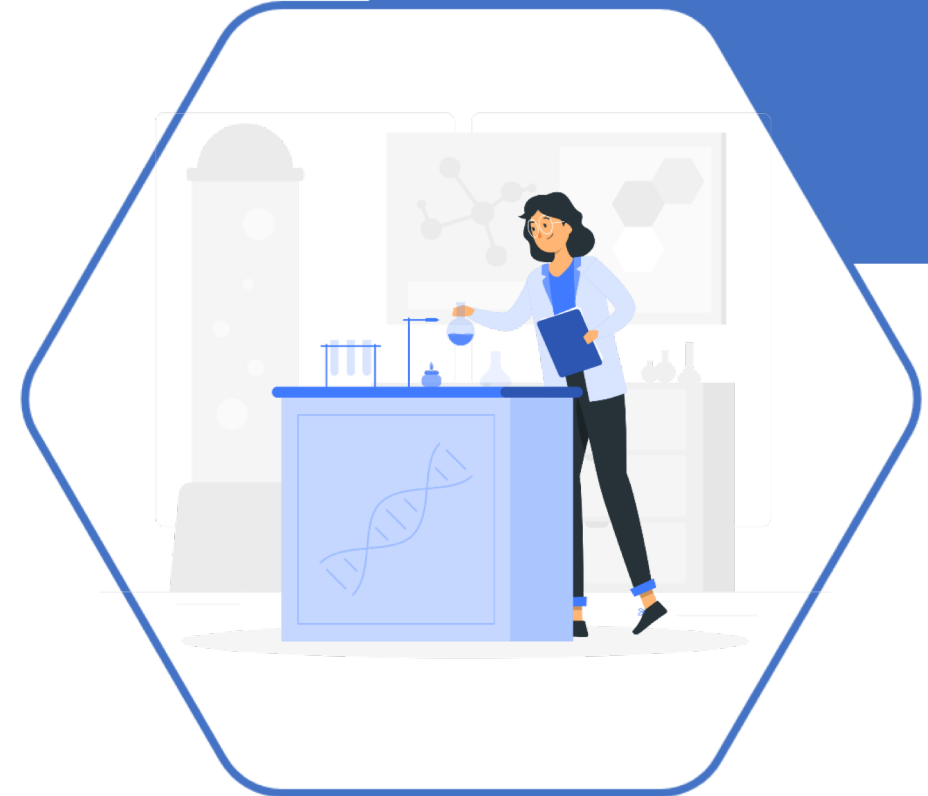


Creative Commons

This work is licensed under a [Creative Commons Attribution-Share Alike Licence \(CC BY-SA 4.0\)](https://creativecommons.org/licenses/by-sa/4.0/).



**Attribution-ShareAlike 4.0 International
(CC BY-SA 4.0)**



Thank you

