



Bacterial strain Taxonomy for Genomic Surveillance

Phylogenetic taxonomies based on whole-genome SNPs

# **Intended Learning Objectives**



Specific objectives of this session:

- 1. Explain the principles of SNP-based taxonomies
- 2. Explore SNP-based nomenclature systems to classify bacterial isolates
- 3. Consider the strengths and limitations of SNP-based approaches in the context of public health microbiology

#### **Outline**



This session consists of the following elements:

- 1. Overview of SNP-based approaches
- 2. Taxonomy and variation between different pathogens
- 3. Inclusion of novel genomes into established nomenclature
- 4. Considering backwards compatibility with previous schemes
- 5. Approaches to define and name novel lineages of concern





# Have you had much experience with SNP-based approaches?



# **Insights from bacterial genomes**

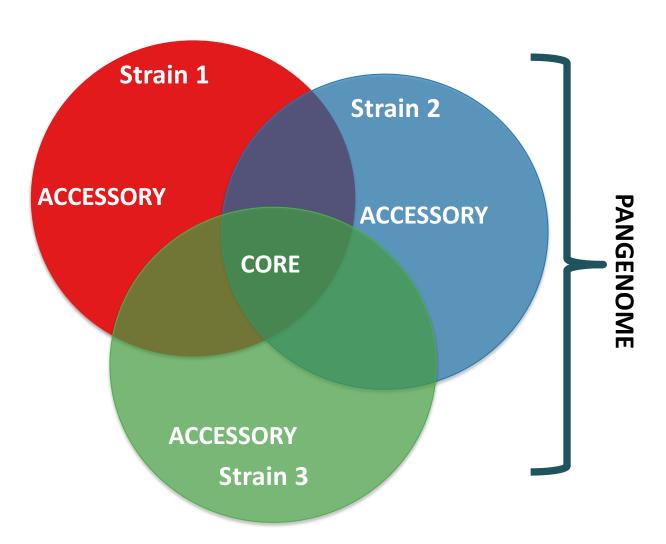


#### **Core genome**

- Present in all isolates
- Infer how related isolates are to each other
- Measurably evolving populations molecular clock of the bacteria
- Vertical evolution

#### **Accessory genome**

- Variable genome content
- Enable adaptation to different ecological niches e.g. AMR genes
- Horizontal evolution



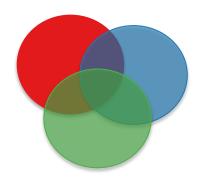
# Genome size differs by bacterial pathogen





#### **Specialist**

- Often have 'closed pangenomes'
- E.g *Coxiella burnetii* genomes ~2,100 genes, obligate intracellular pathogen



#### **Generalists**

- Often have 'open pangenomes'
- E.g. Klebsiella pneumoniae genomes ~5000-6000 genes, multiple ecological niches

Core genome may vary in size depending on bacteria **and** dataset

# Taxonomy differs by bacterial pathogens

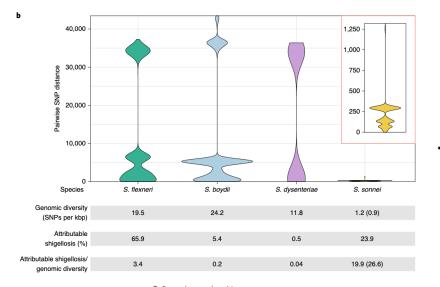


- Important to know the pathogen-specific features of interest
- Shigella
  - 4 species defined by gain of virulence plasmid All sit within *Escherichia coli*

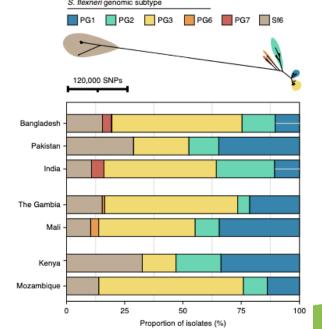
  - Difference in nomenclature within each species
- Klebsiella

  - Big changes to taxonomy over the past decade Multiple new species Importance of clonal groups Use of LIN (Life Identification Number) codes for diverse population from cgMLST
- Salmonella

  - 2 species (discussion if should be more)
    Multiple subspecies within *S. enterica*2,500 serovars defined by surface antigens
    Serovars Typhoidal, non-typhoidal (NTS) and invasive NTS



*S. sonnei* less genetic diversity



Phylogroups within S. flexneri

# Taxonomy differs by bacterial pathogens

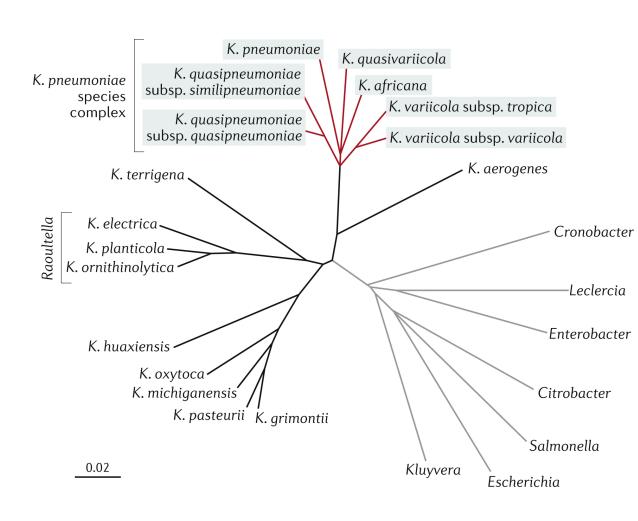


- Important to know the pathogen-specific features of interest
- Shigella
  - 4 species defined by gain of virulence plasmid All sit within *Escherichia coli*

  - Difference in nomenclature within each species
- Klebsiella

  - Big changes to taxonomy over the past decade Multiple new species Importance of clonal groups Use of LIN (Life Identification Number) codes for diverse population from cgMLST
- Salmonella

  - 2 species (discussion if should be more)
    Multiple subspecies within *S. enterica*2,500 serovars defined by surface antigens
    Serovars Typhoidal, non-typhoidal (NTS) and invasive NTS



# Taxonomy differs by bacterial pathogens

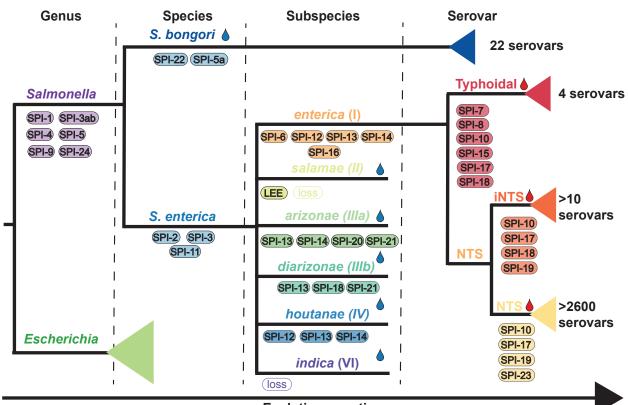


- Important to know the pathogen-specific features of interest
- Shigella
  - 4 species defined by gain of virulence plasmid All sit within *Escherichia coli*

  - Difference in nomenclature within each species
- Klebsiella
  - Big changes to taxonomy over the past decade

  - Multiple new species
    Importance of clonal groups
    Use of LIN (Life Identification Number) codes for diverse population from cgMLST
- Salmonella

  - 2 species (discussion if should be more)
    Multiple subspecies within *S. enterica*2,500 serovars defined by surface antigens
    Serovars Typhoidal, non-typhoidal (NTS) and invasive NTS



**Evolution over time** 

# **Considerations for SNP-based approaches**



- What is important to characterise from the genome?
- How has evolution shaped the bacterial genome (e.g. open / closed)?
- When would SNP-based approaches be useful? Or are there other tools that would suit pathogen better?
- What data would help public health surveillance efforts?
- How should SNP-based approaches be maintained?



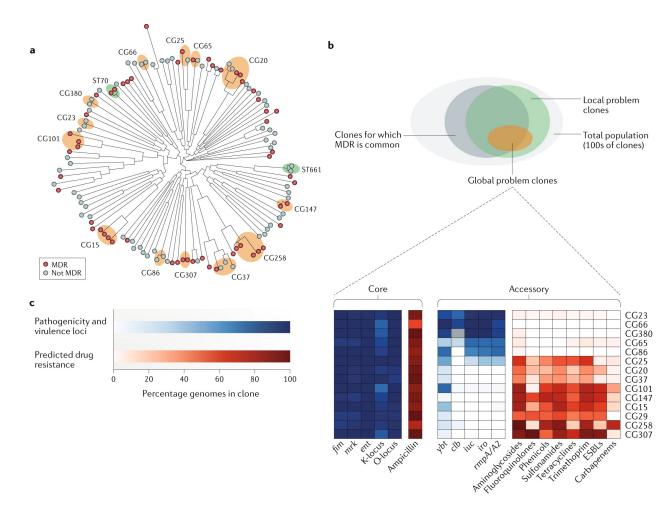


# What are some applications for SNP-based phylogenies?



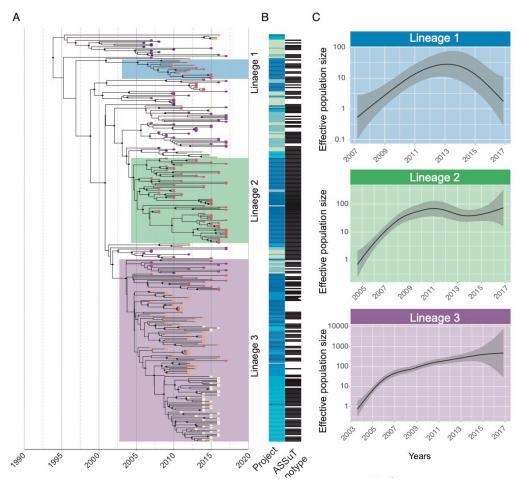
# Different applications of SNP-based phylogenies





#### **Species level - Population framework**

Klebsiella pneumoniae – Clonal Groups (CG) Wyres et al 2020 Nature Reviews Microbiology



#### **Lineage/Serovar Level - Timed tree**

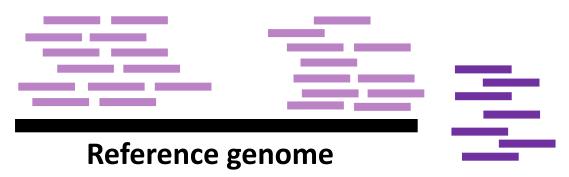
Monophasic *Salmonella* – evolutionary dynamics Ingle et al 2021 *Nature Communications* 

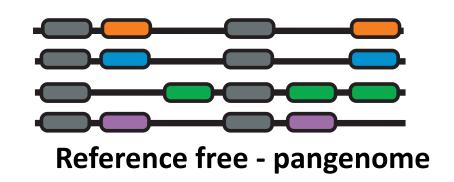
# **SNP-based taxonomies need SNP alignments**



### Considerations for SNP alignments from core genome

- Alignment of reads to reference genome
  - Selection of reference
  - Different tools e.g. Snippy
  - Masking of regions (phage, recombination)
  - Thresholds for filtering of SNPs
- Assembly and pangenome analysis
  - Core genome e.g. of genes in >95% isolates with Panaroo
  - SNPs from core genes





# **How to get SNP-based information**



Alignment of short read data to reference genome

SNPs in the 'core genome' for alignment

Infer (Maximum Likelihood) tree

Identification of clusters/ groups from tree

Use of nomenclature to identify important groups for surveillance

Development of SNP-based methods





# What are some points to consider for SNP-based approaches



# How to get SNP-based information



Alignment of short read data to reference genome

← Selection of reference genome, inclusion of contextual genomes

SNPs in the 'core genome' for alignment

← Size of core genome in dataset

Infer (Maximum Likelihood) tree

← Computational resources available

Approach to take to define e.g. BAPs or SNP threshold→

Identification of clusters/ groups from tree

Development of nomenclature for broader use →

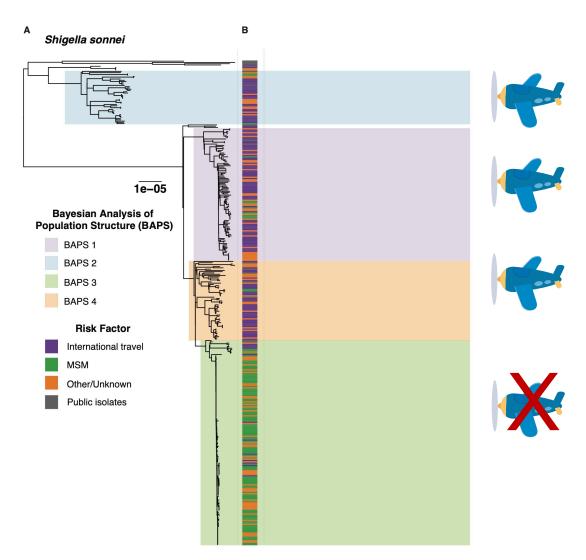
Use of nomenclature to identify important groups for surveillance

Development of SNP-based approaches that don't need a tree every time ->

Development of SNP-based methods

## **Example of SNP-based approaches:** Shigella sonnei



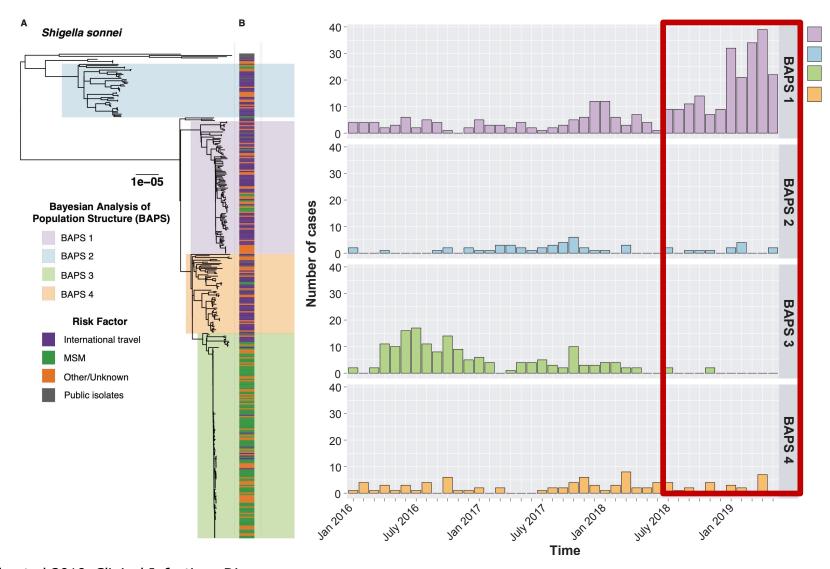


- Increase in cases since 2015
- New epidemiological risk factor
- Infer ML tree and BAPS groups
- Lineages associated with different AMR profiles

Ingle et al 2019 Clinical Infectious Diseases

## **Example of SNP-based approaches:** Shigella sonnei





2019: new year, new outbreak lineage

BAPS 1

BAPS 2

BAPS 3

BAPS 4

- Change in AMR profile
- New isolates added to original dataset
- New SNP alignment and new tree
- 2020: new outbreak lineage

## Example of SNP-based approaches: Shigella sonnei



- Increase in drug-resistant *S. sonnei* globally
- Need approaches to track international spread and standard nomenclature

nature communications



nature communications



Article

https://doi.org/10.1038/s41467-023-37672-w

Article

https://doi.org/10.1038/s41467-023-36222-8

# The evolution and international spread of extensively drug resistant *Shigella sonnei*

Received: 1	2 Septembe	2022

Accepted: 24 March 2023

Published online: 08 April 2023

Check for updates

Lewis C. E. Mason © <sup>1,2</sup>, David R. Greig<sup>3</sup>, Lauren A. Cowley<sup>4</sup>, Sally R. Partridge © <sup>5,6,7,8</sup>, Elena Martinez<sup>7,9</sup>, Grace A. Blackwell<sup>7,9</sup>, Charlotte E. Chong<sup>2</sup>, P. Malaka De Silva<sup>2</sup>, Rebecca J. Bengtsson<sup>2</sup>, Jenny L. Draper © <sup>7,9</sup>, Andrew N. Ginn<sup>7,8,9,10</sup>, Indy Sandaradura © <sup>6,7,9</sup>, Eby M. Sim<sup>5,6,8</sup>, Jonathan R. Iredell <sup>5,6,7,8</sup>, Vitali Sintchenko © <sup>5,6,7,8,9,11</sup>, Danielle J. Ingle © <sup>12</sup>, Benjamin P. Howden © <sup>12</sup>, Sophie Lefèvre © <sup>13</sup>, Elisabeth Njamkepo © <sup>13</sup>, François-Xavier Weill © <sup>13</sup>, Pieter-Jan Ceyssens <sup>14</sup>, Claire Jenkins & Kate S. Baker © <sup>1,2</sup>

#### Aiticie

#### Rapid emergence of extensively drugresistant *Shigella sonnei* in France

Received: 22 November 2022

Accepted: 19 January 2023

Published online: 28 January 2023

Sophie Lefèvre ® ¹, Elisabeth Njamkepo¹, Sarah Feldman ® ².³, Corinne Ruckly¹, Isabelle Carle¹, Monique Lejay-Collin¹, Laëtitia Fabre¹, Iman Yassine ® ¹, Lise Frézal¹, Maria Pardos de la Gandara ® ¹, Arnaud Fontanet² & François-Xavier Weill ® ¹ ⊠

## What information do we want from Typhi genome?



#### Lineage – genotype

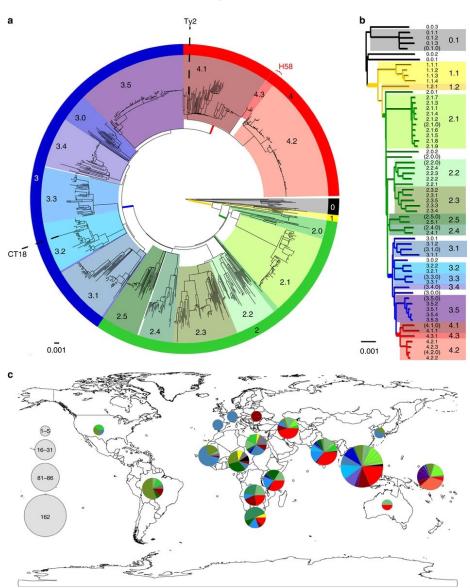
- GenoTyphi scheme
- https://github.com/typhoidgenomics/ genotyphi
- 4 major lineages, > 80 genotypes
- Widely adopted in research and public health settings

#### AMR profile

- SNVs
- Genes

#### Plasmid profile

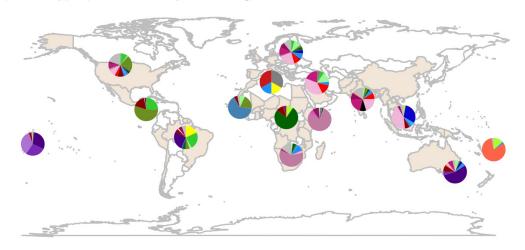
Epidemiologically important plasmids



# **Insights from 13,000 Typhi genomes**



#### a) Genotype prevalence by world region, 2010-2020



Collaborative efforts >200 researchers/ public health workers

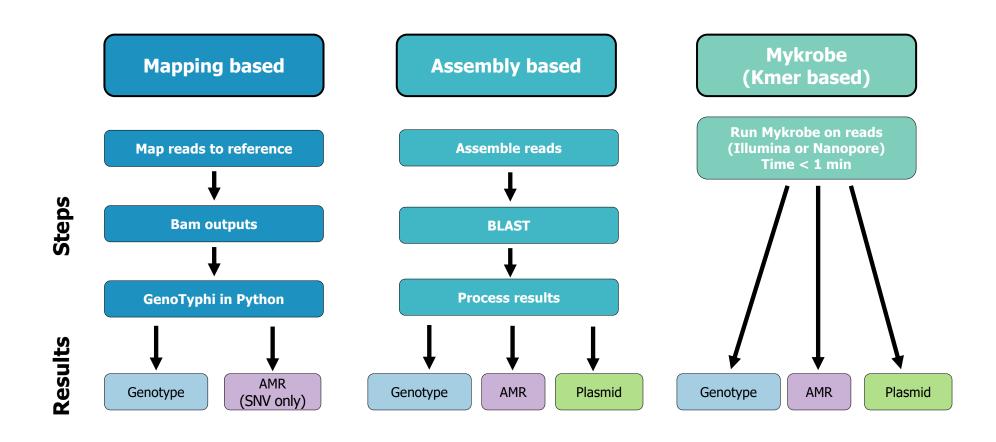
Enhanced genomic surveillance

a) Annual genotype prevalence, by lab Bangladesh Bangladesh Bangladesh Bangladesh Bangladesh Bangladesh India CDC CHR **CMH** MDU ICD PHE 0.75 -0.50 -0.25 -0.00 -India India India India India India India **BCH** CDC CMC CML CNH **ESR** KHM 1.00 -0.75 -0.50 -0.25 -0.00 -India India India India India India India MDU PHE **KKC** PGI RDT SAP SJB 0.75 -0.50 -0.25 -0.00 -Nepal Pakistan Pakistan **Pakistan Pakistan Pakistan** Nepal **KUH** PAH AKU CDC MDU PHE 1.00 **-**0.75 **-**

Carey et al. 2023 *eLife*https://www.typhoidgenomics.org/
https://github.com/typhoidgenomics/TyphoidGenomicsConsortiumWG1

# How to get SNP-based taxonomy information

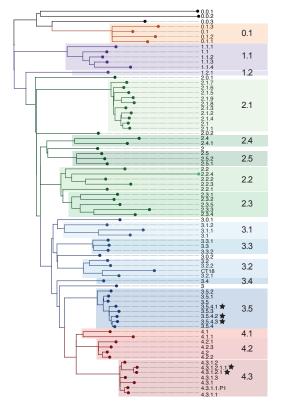




Genotyping schemes can be implemented on different platforms

# **Typhi Mykrobe**

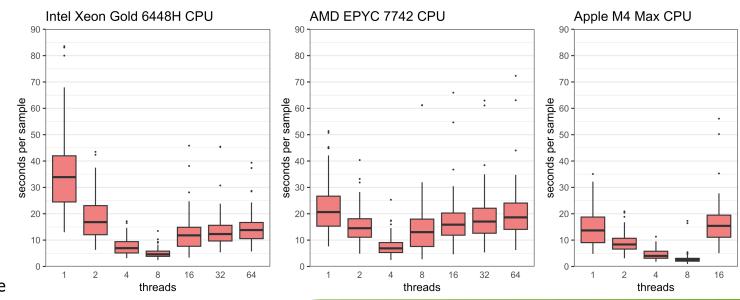
- Implementation of GenoTyphi scheme in Mykrobe
- SNP-based approach use of established poplation framework
- Fast. Typically <1 min genome
- Enables inference of tree position without the inferring a tree each time
- Provides standard nomenclature for lineages and sub- lineages





#### **GenoTyphi framework**

- 87 genotypes
- 4 primary lineages
- 16 clades
- 49 subclades



# **Running Typhi Mykrobe**



# Installing mykrobe First, install Mykrobe (v0.10.0+) as per the instructions on the Mykrobe github. Once Mykrobe is installed, you can run the following two commands to ensure you have the most up-to-date panels for genotyping, including the Typhi panel (latest version, v20240407): mykrobe panels update\_metadata mykrobe panels update\_species all You can check what version of the scheme is currently loaded in your Mykrobe installation via:

```
Run Mykrobe on fastq file/s for a given genome:

mykrobe predict --sample aSample \
--species typhi \
--format json \
--out aSample.json \
--seq aSample_1.fastq.gz aSample_2.fastq.gz
```

```
Tabulate Mykrobe results for one or more genomes:

(requires Python3 + pandas library)

(python script parse_typhi_mykrobe.py is in this repository in the /typhimykrobe directory in this repository)

python parse_typhi_mykrobe.py --jsons *.json --prefix mykrobe_out
```

mykrobe panels describe

# **Typhi Mykrobe**



Direct from reads (Illumina or ONT)

#### Output – tabular format

#### Mykrobe parse typhi mykrobe.py invA present (to confirm Salmonella enterica) Summarise and tabulate → JSON Species and serovar calls S. enterica MLST match (to confirm Typhi) **JSON** Genotype calls reads .fastq Genotype confidence (input files) Type GenoTyphi marker (calculate hierarchical genotype) → JSON AMR determinants by drug Plasmid markers **Detect AMR genes and SNVs** (organise by drug) (one per genome) **Detect plasmid replicons + pST6** (report inidividually) Results table (TSV)

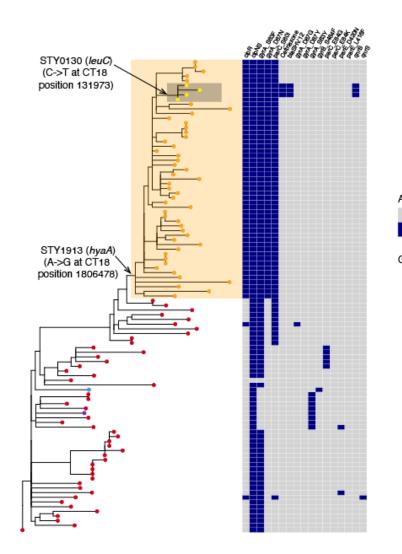
				lowest support for	supported	for	additional	
species	spp_percent	genotype	confidence	genotype marker	markers	additional	markers	node support
typhi	91.518	4.3.1.1	strong	-	-	-	-	1 (1; 0/159); 2 (1; 0/180); 3 (1; 0/155); 4 (1; 133/0); 4.3.1 (1; 138/0); 4.3.
typhi	91.565	4.3.1.2	strong	-	-	-	-	1 (1; 0/167); 2 (1; 0/176); 3 (1; 0/146); 4 (1; 155/0); 4.3.1 (1; 153/0); 4.3.
typhi	91.723	2	2 strong	-	-	-	-	1 (1; 0/133); 2 (1; 1/138); 2.2 (1; 129/0)
typhi	91.504	4.3.1.2	strong	-	-	-	-	1 (1; 0/133); 2 (1; 0/147); 3 (1; 0/160); 4 (1; 139/0); 4.3.1 (1; 121/0); 4.3.
typhi	91.48	4	1 strong	-	-	-	-	1 (1; 0/148); 2 (1; 0/137); 3 (1; 0/152); 4 (1; 164/0); 4.1 (1; 120/0)
typhi	91.892		3 strong	-	-	-	-	1 (1; 0/61); 2 (1; 0/63); 3 (1; 0/75)
typhi	91.478		3 strong	-	-	-	-	1 (1; 0/57); 2 (1; 0/62); 3 (1; 0/66)
typhi	91.494	4	1 strong	-	-	-	-	1 (1; 0/86); 2 (1; 1/58); 3 (1; 0/88); 4 (1; 48/1); 4.1 (1; 56/0)
typhi	91.746	3.2.1	strong	-	-	-	-	1 (1; 1/63); 2 (1; 0/63); 3 (1; 0/63); 3.2 (1; 0/54); 3.2.1 (1; 0/67)
typhi	91.715	3.2.1	strong	-	-	-	-	1 (1; 0/71); 2 (1; 0/58); 3 (1; 0/80); 3.2 (1; 1/57); 3.2.1 (1; 0/62)

A. Typhi Mykrobe pipeline

# **Ongoing efforts for genotyping Typhi**



- Typhi Mykrobe led by a working group within Global Typhoid Genomics Consortium
- Another working group exploring how to add new genotypes
  - GenoTýphi Genotyping Scheme Development and Curation Working Group
- Considerations for incorporating new genomes into the taxonomy and genotypes
   Need to new tree to identify SNPs specific
  - Need to new tree to identify SNPs specific to lineages
- Considerations for how to name lineages going forward



# SNP-based schemes useful for clonal pathogens



#### Mykrobe schemes already exist for:

- Mycobacterium tuberculosis
- Staphylococcus aureus
- Shigella sonnei
- Salmonella enterica serovar Paratyphi B
- Salmonella enterica serovar Typhi

# Mykrobe schemes in development for:

• Shigella flexneri

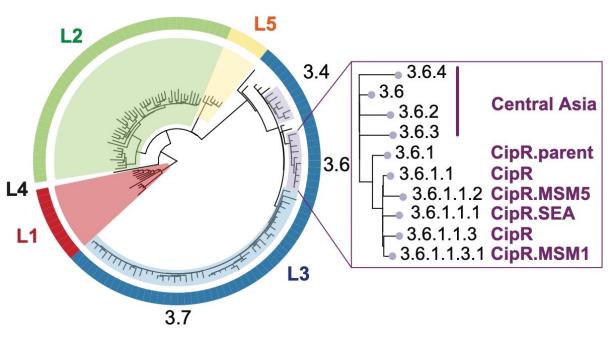
#### Considerations for developing schemes

- Need discovery tree to define genotypes and identify SNPs
- Use of previously identified lineages/ clusters to inform design
- Naming of lineages
- When to define new lineages

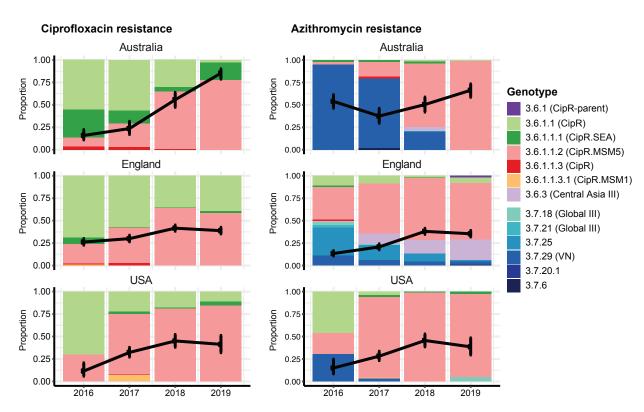
# Development of genotype schemes implemented in Mykrobe



- Genotyping scheme developed for Shigella sonnei
- Suited as clonal pathogen



Five lineages with 128 genotypes



Applied to >4,000 genomes from 3 public health labs

28





# What should be considered when naming a novel lineage?





# In summary



List of learning points in this session:

- SNP-based approaches can be useful in tracking priority bacterial pathogens
  - Trees remain important for defining SNP-based genotyping schemes
  - Genotyping schemes implemented in Mykrobe can provide rapid and robust lineage identification for public health surveillance
- Taxonomy and nomenclature varies between bacterial pathogens
  - This is underpinned by evolution of the pathogen
  - Important to consider the dataset and research/ public health questions when using SNP-based approaches

# **Further reading**



- https://github.com/rrwick/Core-SNP-filter
- https://github.com/gtonkinhill/rhierbaps
- https://github.com/gtonkinhill/panaroo
- https://github.com/Mykrobe-tools/mykrobe
- https://github.com/typhoidgenomics/genotyphi
- https://github.com/typhoidgenomics/TyphoidGenomicsConsortiumMykrobe
- https://github.com/katholt/sonneityping

#### References



- Wong et al. 2016. An extended genotyping framework for Salmonella enterica serovar Typhi, the cause of human typhoid. Nature Communications. <a href="https://doi.org/10.1038/ncomms12827">https://doi.org/10.1038/ncomms12827</a>
- Ingle et al. Typhi Mykrobe: fast and accurate lineage identification and antimicrobial resistance genotyping directly from sequence reads for the typhoid fever agent Salmonella Typhi. BioRxiv: <a href="https://doi.org/10.1101/2024.09.30.613582">https://doi.org/10.1101/2024.09.30.613582</a>
- Ingle et al. 2019. Co-circulation of Multidrug-resistant Shigella Among Men Who Have Sex With Men in Australia. Clinical Infectious Diseases. <a href="https://doi.org/10.1093/cid/ciz005">https://doi.org/10.1093/cid/ciz005</a>
- Hawkey et al. 2021. Global population structure and genotyping framework for genomic surveillance of the major dysentery pathogen, Shigella sonnei. *Nature Communications*. <a href="https://doi.org/10.1038/s41467-021-22700-4">https://doi.org/10.1038/s41467-021-22700-4</a>







# Acknowledgements

The creation of this training material was commissioned by ECDC to The University of Melbourne and the Peter Doherty Institute for Infection and Immunity with the direct involvement of Dr. Danielle Ingle