



Bacterial strain Taxonomy for Genomic Surveillance

## ReporTree: Clustering and Visualization for Surveillance (practical)

## **Intended Learning Objectives**



Specific objectives of this session:

- 1. Learn what is ReporTree and its utility in routine genomic surveillance
- 2. Learn the principles behind ReporTree cluster nomenclature
- 3. Learn how clustering stability regions can be used for nomenclature design
- 4. Learn how to launch a ReporTree command line for routine surveillance

#### **Outline**



This session consists of the following elements:

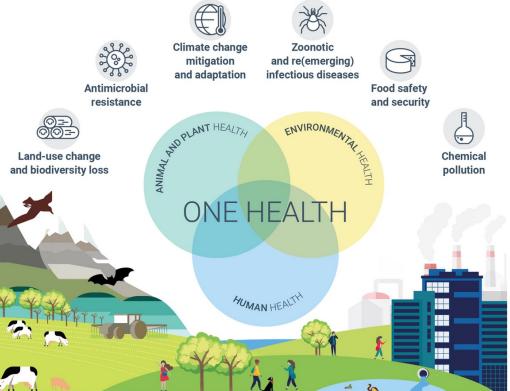
- 1. Introduction to ReporTree
- 2. ReporTree's input flexibility and main outputs
- 3. Detect and track genetic clusters with ReporTree nomenclature system
- 4. Identification of candidate thresholds for nomenclature design (clustering stability regions)
- 5. Exercise of the real-life application of ReporTree in routine surveillance

#### **Genomic surveillance**



Genomic surveillance is the process of constantly monitoring pathogens and analyzing their genetic similarities and differences.

WHO, 2022

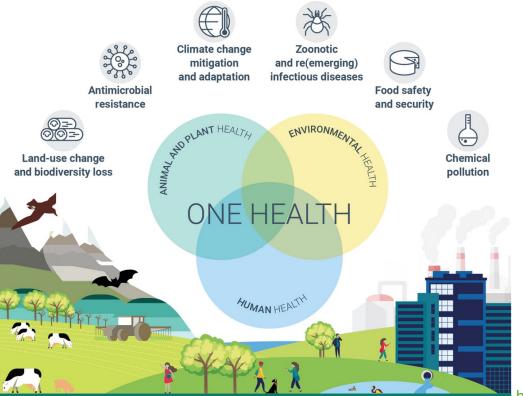


#### **Genomic surveillance**



Genomic surveillance is the process of constantly monitoring pathogens and analyzing their genetic similarities and differences.

WHO, 2022



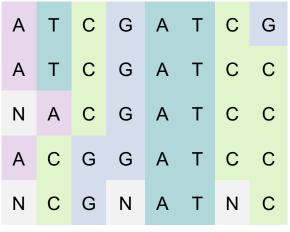
A pivotal outcome of genomics surveillance is the identification of pathogen genetic clusters/lineages and their characterization in terms of geotemporal spread or linkage to clinical and demographic data

## **Analyses of Whole-Genome Sequencing data**



Multiple bioinformatics solutions, but similar steps and goals...

#### **Genetic data**



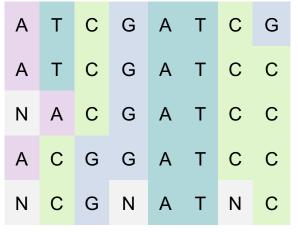
SNPs\* or alleles

## **Analyses of Whole-Genome Sequencing data**



Multiple bioinformatics solutions, but similar steps and goals...

#### **Genetic data**



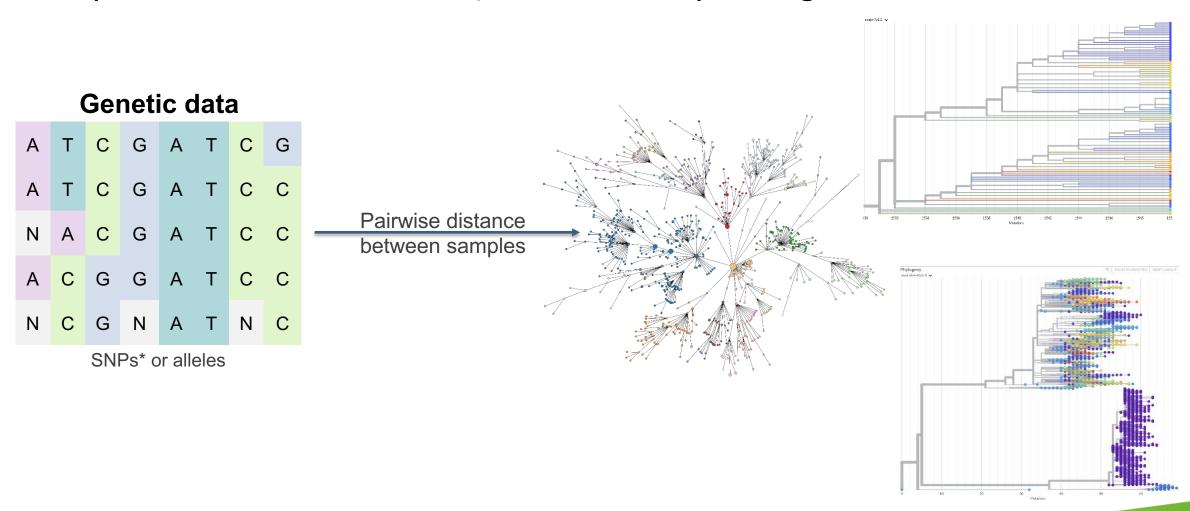
SNPs\* or alleles

Pairwise distance between samples

## **Analyses of Whole-Genome Sequencing data**



Multiple bioinformatics solutions, but similar steps and goals...







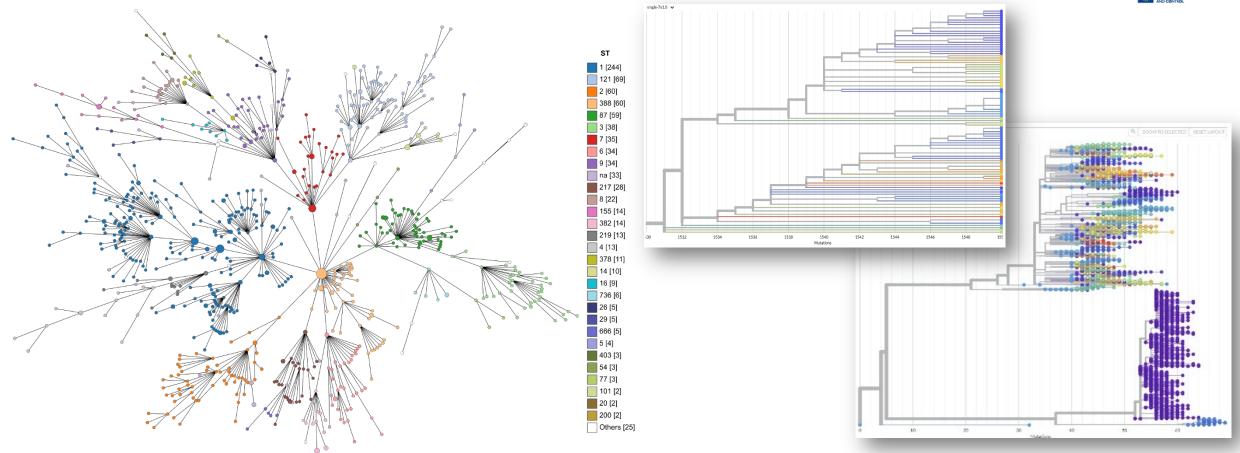
# How do you perform bacterial clustering analyses?





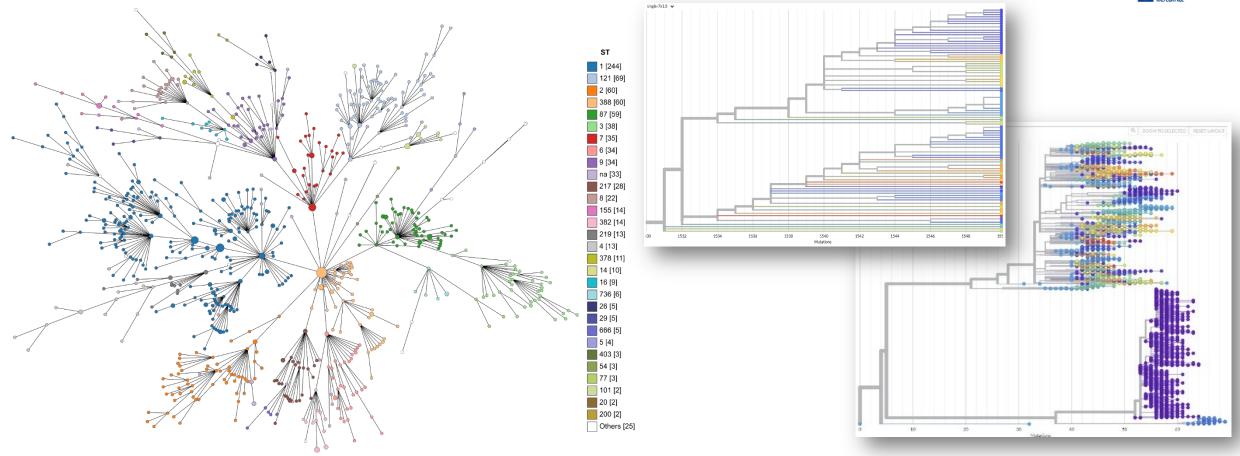
## Detection of genetic clusters and characterization exc





## Detection of genetic clusters and characterization





The detection of genetic clusters and their characterization with clinical/epidemiological data is still a challenging step that often relies on non-automated approaches

#### ReporTree



- Flexible solution to <u>automatically identify genetic clusters</u> at any (or all) distance thresholds
- <u>Generate surveillance-oriented reports</u> based on the available metadata, such as timespan, geography or clinical status



#### ReporTree



- Flexible solution to <u>automatically identify genetic clusters</u> at any (or all) distance thresholds
- Generate surveillance-oriented reports based on the available metadata, such as timespan, geography or clinical status





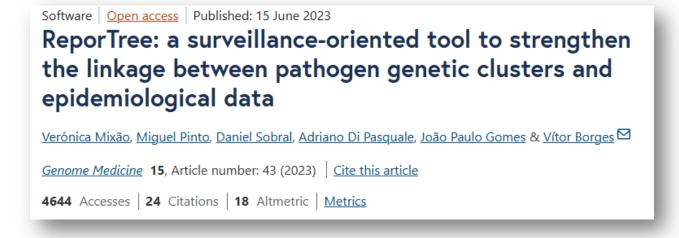


#### ReporTree



- Flexible solution to <u>automatically identify genetic clusters</u> at any (or all) distance thresholds
- Generate surveillance-oriented reports based on the available metadata, such as timespan, geography or clinical status







Also available through: PubMLST, COHESIVE, etc.





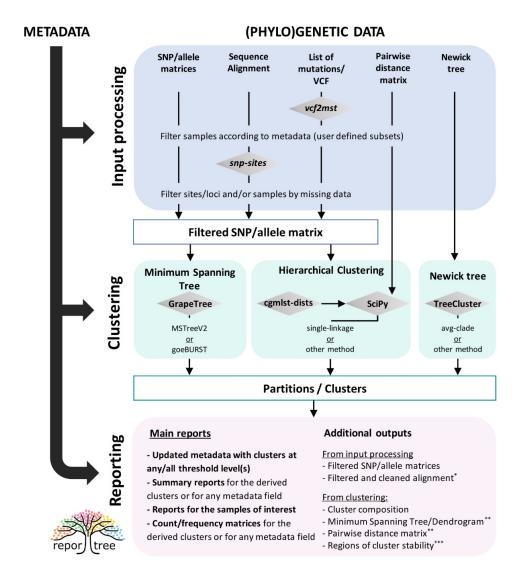
# Which WGS-based typing approach do you use for routine surveillance of bacterial pathogens?





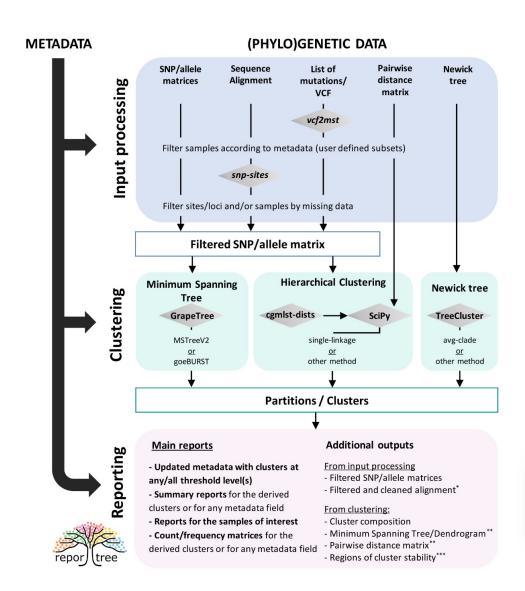
## Designed for surveillance of multiple pathogens





## Designed for surveillance of multiple pathogens





#### Main outputs:

- Tree
- Metadata table updated with clusters at any/all resolution levels
- Summary reports with the statistics/trends for the derived genetic clusters
- Summary reports for the samples of interest (e.g., new samples)

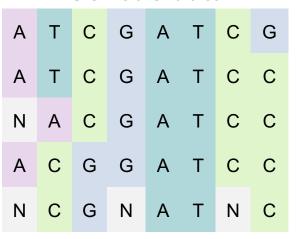
partition	cluster	cluster_ length	samples	source	country	first_seq_ date	last_seq_ date
MST-4x1.0	cluster_34	2	sample_0419,sample_0464	clinical (50.0%), food (50.0%)	C (50.0%), B (50.0%) (n = 2)	20/06/11	12/03/12
MST-7x1.0	cluster_106	33	sample_0017,sample_0037,sample_	clinical (63.6%), food (36.4%)	C (42.4%), B (30.3%), A (27.3%)	30/03/04	14/07/21
MST-14x1.0	cluster_87	46	sample_0017,sample_0037,sample_	clinical (67.4%), food (32.6%)	A (37.0%), C (34.8%), B (28.3%)	30/03/04	14/07/21



## How does it work?



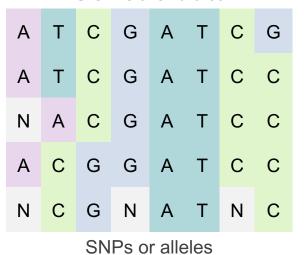
#### **Genetic data**



SNPs or alleles



#### **Genetic data**



Treated similar to wgMLST data (alignment\_processing.py)



#### **Genetic data**

Α	Т	С	G	Α	Т	С	G
Α	Т	С	G	Α	Т	С	С
N	Α	С	G	Α	Т	С	С
Α	С	G	G	Α	Т	С	С
N	С	G	N	Α	Т	N	С

SNPs or alleles

cgMLST (set of fixed core loci)

wgMLST (set of core loci + accessory loci)

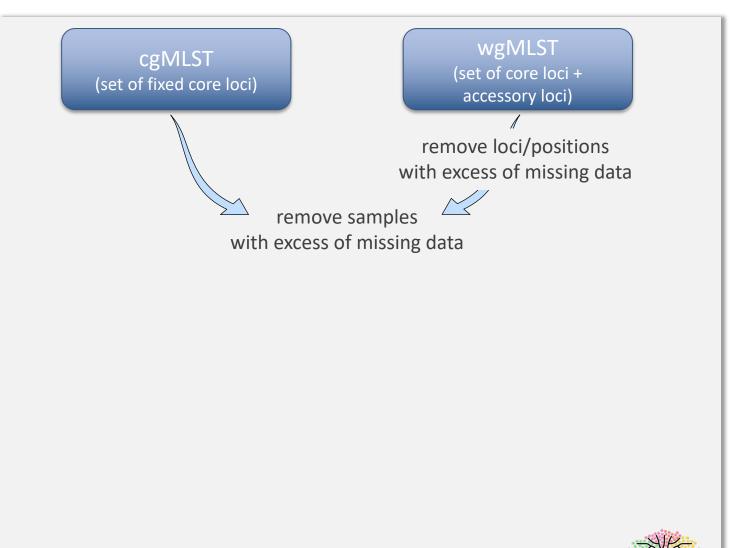




#### **Genetic data**

Α	Т	С	G	Α	Т	С	G
Α	Т	С	G	Α	Т	С	С
N	Α	С	G	Α	Т	С	С
Α	С	G	G	Α	Т	С	С
N	С	G	N	Α	Т	N	С

SNPs or alleles

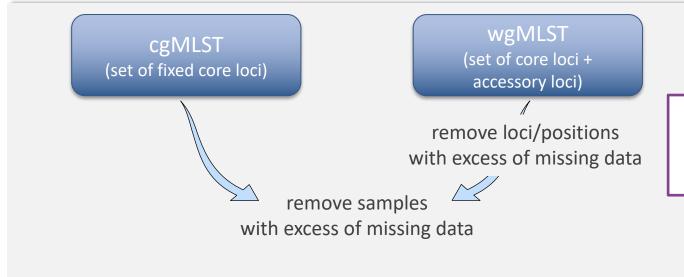




#### **Genetic data**

Α	Т	С	G	Α	Т	С	G
Α	Т	С	G	Α	Т	С	С
N	Α	С	G	Α	Т	С	С
Α	С	G	G	Α	Т	С	С
N	С	G	N	Α	Т	N	С

SNPs or alleles



or based on a list of loci that represent a static cgMLST schema!

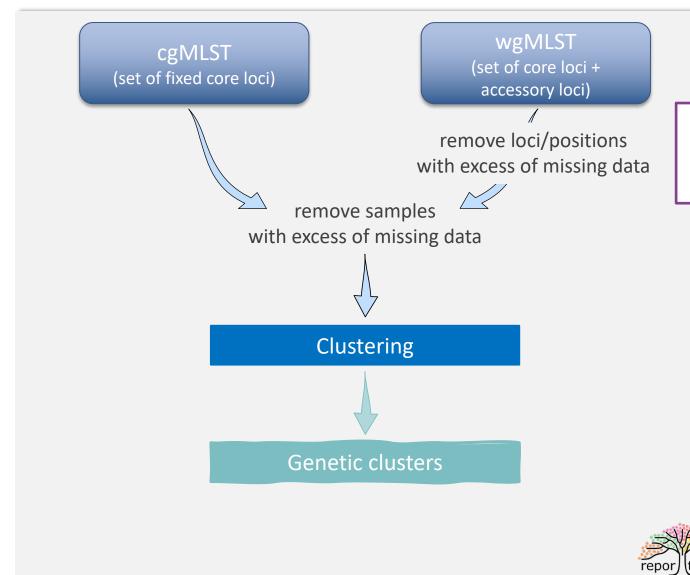




#### **Genetic data**

Α	Т	С	G	Α	Т	С	G
Α	Т	С	G	Α	Т	С	С
N	Α	С	G	Α	Т	С	С
Α	С	G	G	Α	Т	С	С
N	С	G	N	Α	Т	N	С

SNPs or alleles



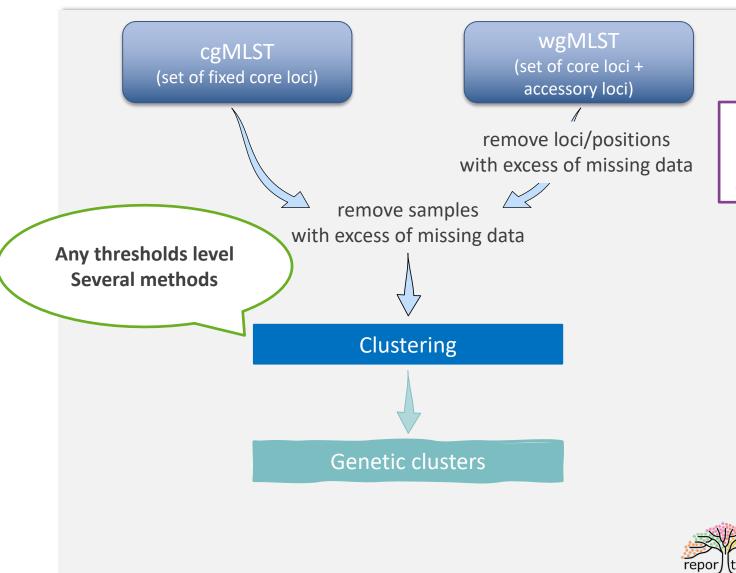
or based on a list of loci that represent a static cgMLST schema!



#### **Genetic data**

A T C G A T C G
A T C G A T C C
N A C G A T C C
A C G A T C C
N C G N A T C C

SNPs or alleles



or based on a list of loci that represent a static cgMLST schema!

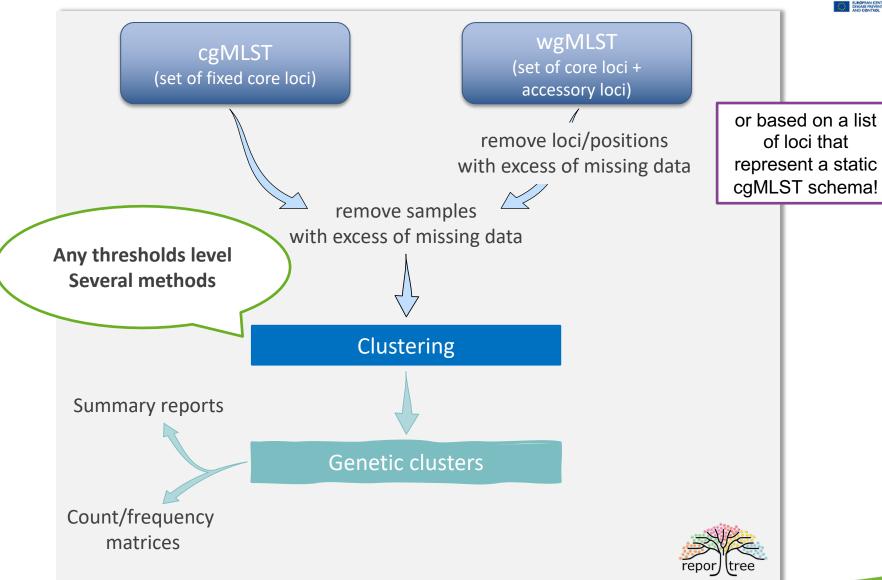


of loci that

#### **Genetic data**

A T C G A T C G T C G A T C C A C G A T C C

SNPs or alleles

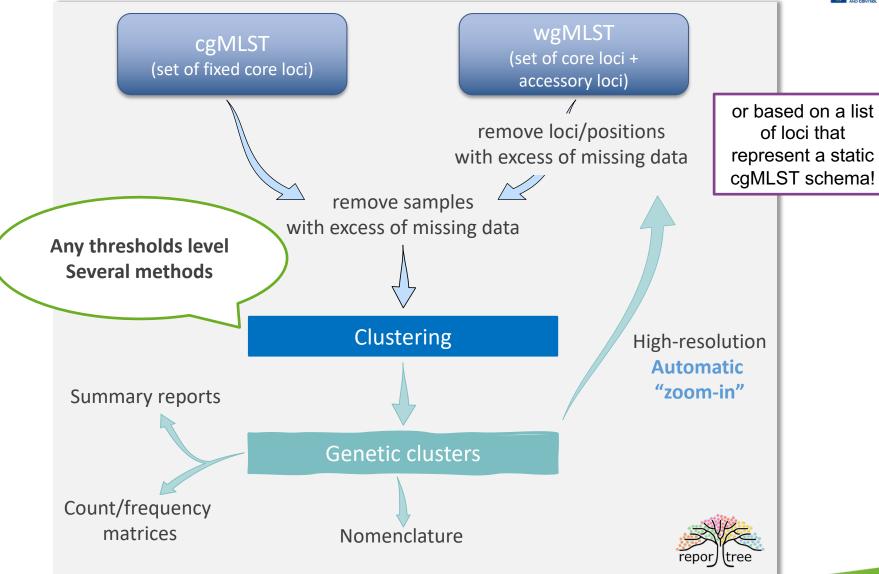




#### **Genetic data**

A T C G A T C G
A T C G A T C C
N A C G A T C C
N C G A T C C
N C G A T C C

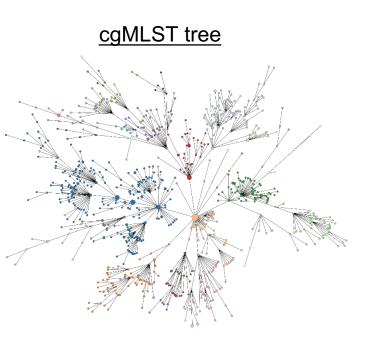
SNPs or alleles



## Dynamic wgMLST approach



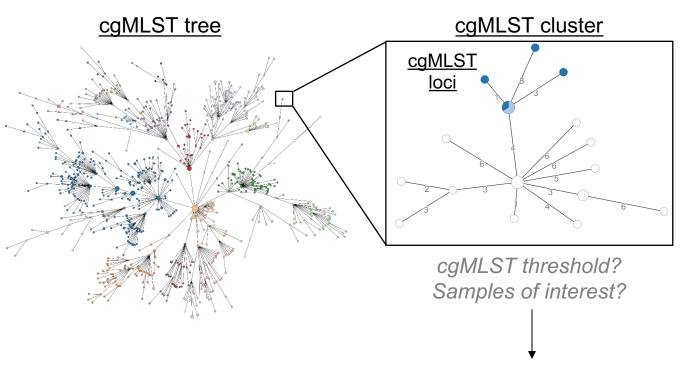
• Useful for species with wgMLST schemas (e.g. Salmonella enterica or Escherichia coli)



#### Dynamic wgMLST approach



• Useful for species with wgMLST schemas (e.g. Salmonella enterica or Escherichia coli)



E.g. zoom-in for all clusters detected at X allelic differences where the samples of this week belong

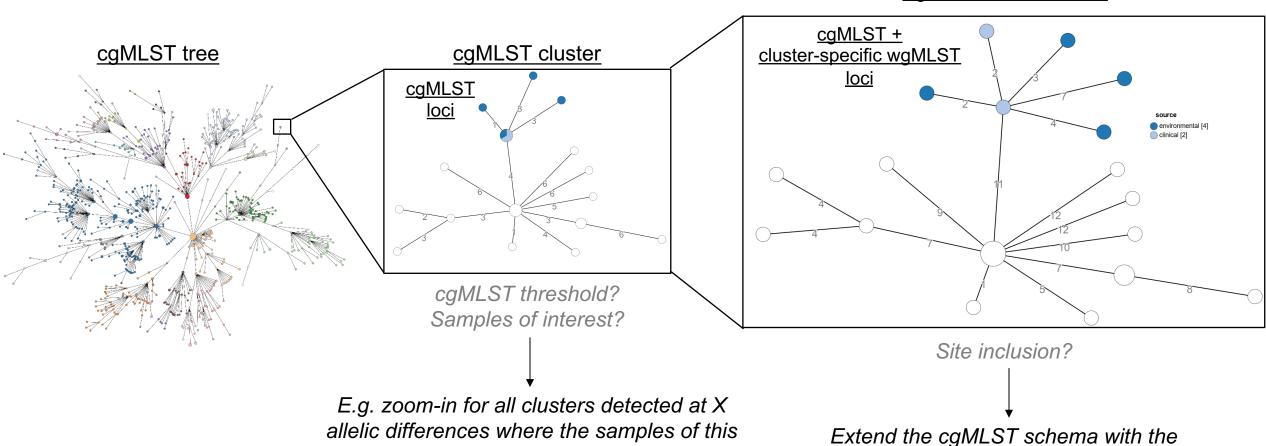
### Dynamic wgMLST approach



• Useful for species with wgMLST schemas (e.g. Salmonella enterica or Escherichia coli)

#### wgMLST cluster zoom-in

accessory loci present in X% of the cluster samples



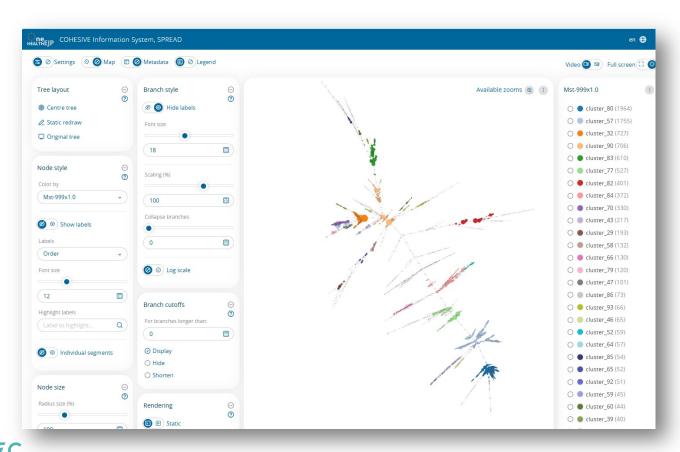
week belong

Mixão et al. 2023 *Genome Med* (doi: 10.1186/s13073-023-01196-1) \*Concept similar to PHYLOViZ (doi: 10.1186/1471-2105-13-87)

### Dynamic wgMLST approach - visualization



 SPREAD - Spatiotemporal Pathogen Relationships and Epidemiological Analysis Dashboard

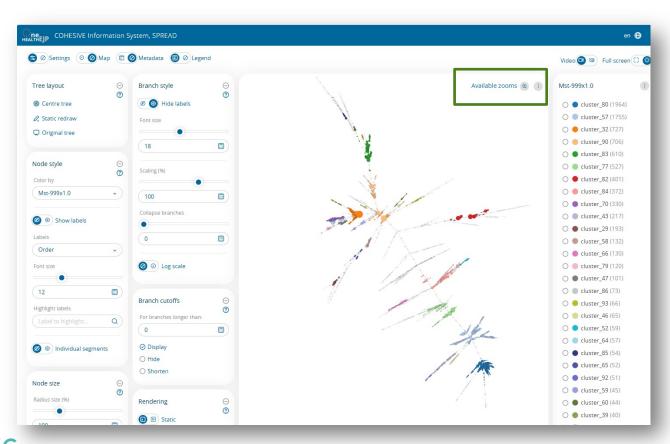




### Dynamic wgMLST approach - visualization



 SPREAD - Spatiotemporal Pathogen Relationships and Epidemiological Analysis Dashboard

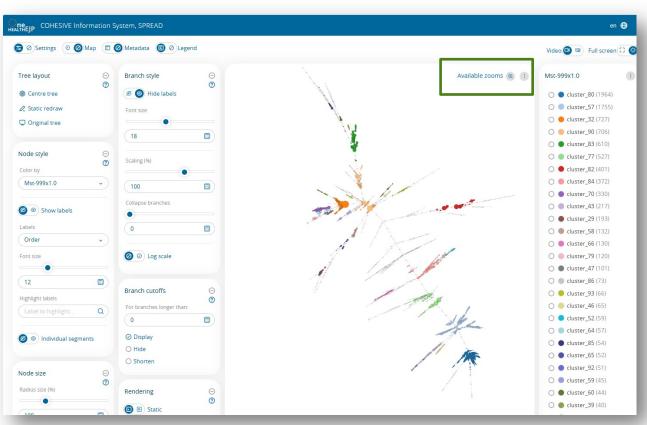


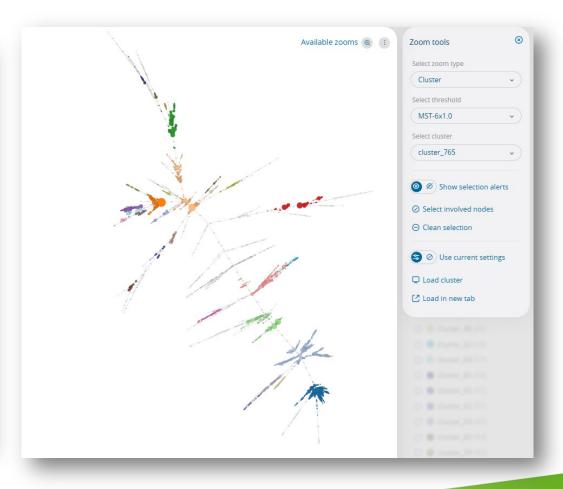


### Dynamic wgMLST approach - visualization



 SPREAD - Spatiotemporal Pathogen Relationships and Epidemiological Analysis Dashboard





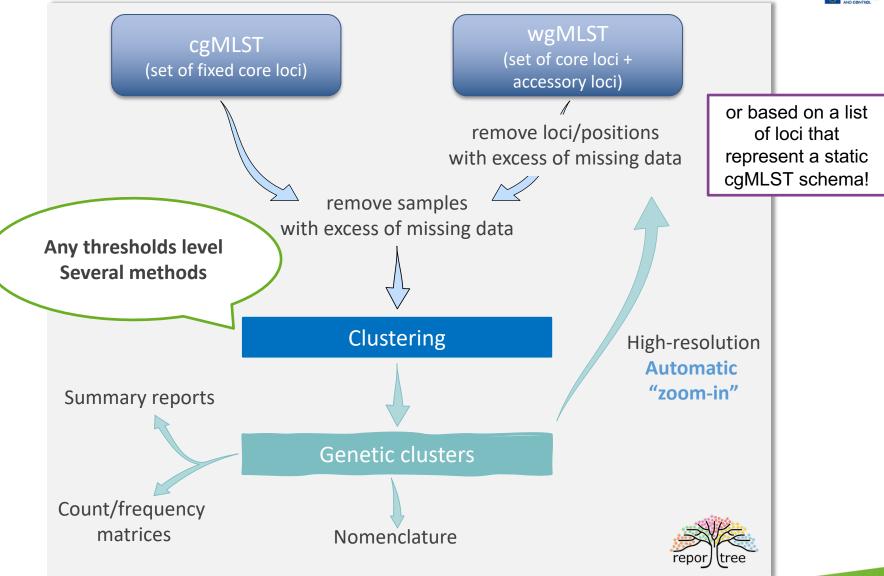
https://github.com/genpat-it/spread



#### **Genetic data**

A T C G A T C G
A T C G A T C C
N A C G A T C C
A C G A T C C
N C G N A T C C

SNPs or alleles

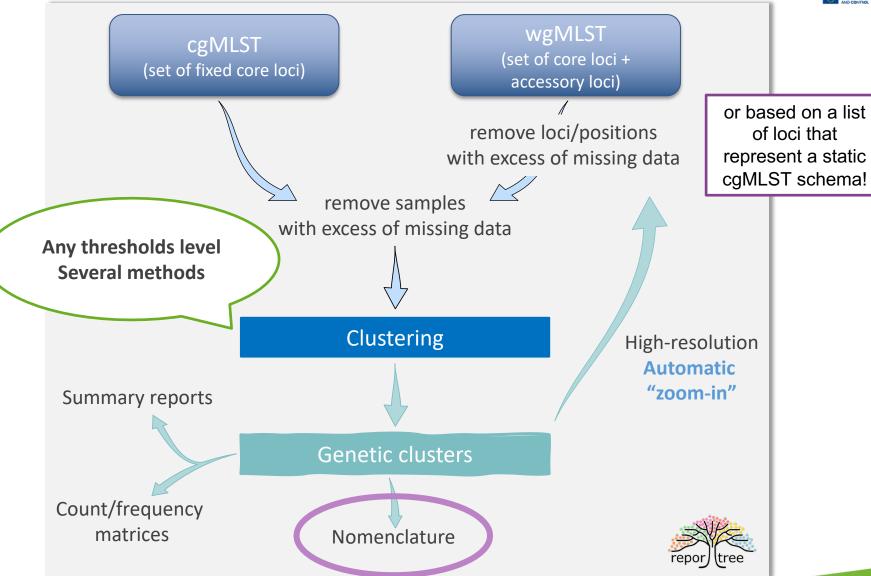




#### **Genetic data**

A T C G A T C G
A T C G A T C C
N A C G A T C C
A C G A T C C
N C G N A T C C

SNPs or alleles





# How does it help me monitoring clusters and communicating?

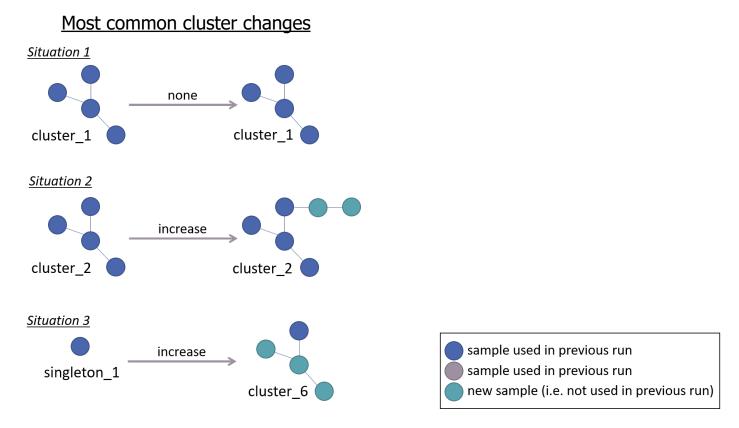


<u>Cluster names are maintained in subsequent runs</u>, facilitating the establishment of a stable local nomenclature...



Cluster names are maintained in subsequent runs, facilitating the establishment of

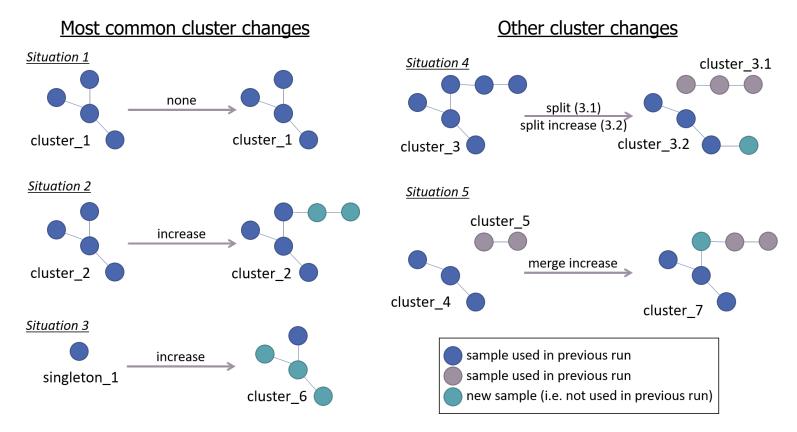
a stable local nomenclature...





<u>Cluster names are maintained in subsequent runs</u>, facilitating the establishment of

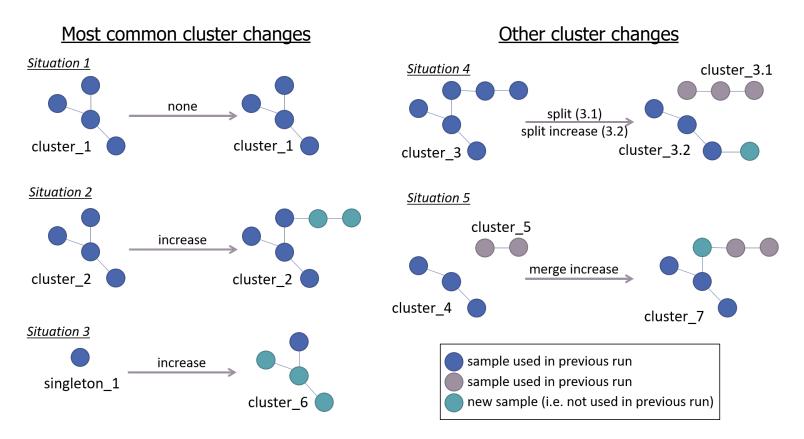
a stable local nomenclature...





<u>Cluster names are maintained in subsequent runs</u>, facilitating the establishment of

a stable local nomenclature...



Cluster names
always reflect the
structure of the tree



<u>Cluster names are maintained in subsequent runs</u>, facilitating the establishment of

a stable local nomenclature...

partition	cluster	nomenclature_change	n_increase	cluster_length	first_seq_date	last_seq_date	source	samples_increase
MST-10x1.0	cluster_1	kept (none)	0	50	14/01/2023	02/06/2025	Clinical (96.0%), Food (4.0%) (n = 50)	
MST-10x1.0	cluster_2	kept (increase)	1	21	07/01/2024	01/09/2025	Clinical (95.2%), Food (4.8%) (n =21)	Sample_A
MST-10x1.0	cluster_3	kept (increase)	2	20	21/02/2025	01/09/2025	Clinical (100.0%) (n = 20)	Sample_B, Sample_C
MST-10x1.0	cluster_4	kept (none)	0	10	07/09/2024	24/06/2025	Clinical (100.0%) (n = 10)	
MST-10x1.0	cluster_5	new	4	4	09/08/2025	01/09/2025	Clinical (100.0%) (n = 4)	Sample_D, Sample_E, Sample_F, Sample_G



<u>Cluster names are maintained in subsequent runs</u>, facilitating the establishment of

a stable local nomenclature...

partition	cluster	nomenclature_change	n_increase	cluster_length	first_seq_date	last_seq_date	source	samples_increase
MST-10x1.0	cluster_1	kept (none)	0	50	14/01/2023	02/06/2025	Clinical (96.0%), Food (4.0%) (n = 50)	
MST-10x1.0	cluster_2	kept (increase)	1	21	07/01/2024	01/09/2025	Clinical (95.2%), Food (4.8%) (n =21)	Sample_A
MST-10x1.0	cluster_3	kept (increase)	2	20	21/02/2025	01/09/2025	Clinical (100.0%) (n = 20)	Sample_B, Sample_C
MST-10x1.0	cluster_4	kept (none)	0	10	07/09/2024	24/06/2025	Clinical (100.0%) (n = 10)	
MST-10x1.0	cluster_5	new	4	4	09/08/2025	01/09/2025	Clinical (100.0%) (n = 4)	Sample_D, Sample_E, Sample_F, Sample_G



<u>Cluster names are maintained in subsequent runs</u>, facilitating the establishment of

a stable local nomenclature...

partition	cluster	nomenclature_change	n_increase	cluster_length	first_seq_date	last_seq_date	source	samples_increase
MST-10x1.0	cluster_1	kept (none)	0	50	14/01/2023	02/06/2025	Clinical (96.0%), Food (4.0%) (n = 50)	
MST-10x1.0	cluster_2	kept (increase)	1	21	07/01/2024	01/09/2025	Clinical (95.2%), Food (4.8%) (n =21)	Sample_A
MST-10x1.0	cluster_3	kept (increase)	2	20	21/02/2025	01/09/2025	Clinical (100.0%) (n = 20)	Sample_B, Sample_C
MST-10x1.0	cluster_4	kept (none)	0	10	07/09/2024	24/06/2025	Clinical (100.0%) (n = 10)	
MST-10x1.0	cluster_5	new	4	4	09/08/2025	01/09/2025	Clinical (100.0%) (n = 4)	Sample_D, Sample_E, Sample_F, Sample_G

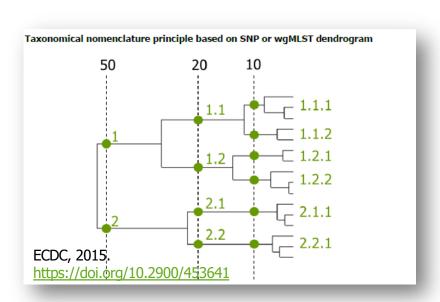
All changes in the tree are recorded in a comprehensive report!



<u>Cluster names are maintained in subsequent runs</u>, facilitating the establishment of a stable local nomenclature... and, if requested, a hierarchical nomenclature similar to "SNP-address".



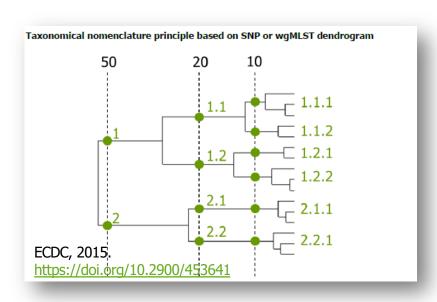
<u>Cluster names are maintained in subsequent runs</u>, facilitating the establishment of a stable local nomenclature... and, if requested, a hierarchical nomenclature similar to "SNP-address".





<u>Cluster names are maintained in subsequent runs</u>, facilitating the establishment of

a stable local nomenclature... and, if requested, a hierarchical nomenclature similar to "SNP-address".



The user can set the hierarchical thresholds for nomenclature (e.g. 50,20,7,0)

sample	50	20	7	0	nomenclature_code_2025-10-01
sample_01	C564	C550	C647	C38	C564-C550-C647-C38
sample_02	C564	C550	C647	C38	C564-C550-C647-C38
sample_03	C564	C550	C647	S1337	C564-C550-C647-S1337
sample_04	S136	S492	S1397	S4880	S136-S492-S1397-S4880
sample_05	C452	S28	S927	S4408	C452-S28-S927-S4408
sample_06	C286	C372	C626	S2861	C286-C372-C626-S2861
sample_07	C286	C372	C626	S2862	C286-C372-C626-S2862
sample_08	C21	C37	C75	C570	C21-C37-C75-C570
sample_09	C21	C37	C75	C570	C21-C37-C75-C570

C - cluster

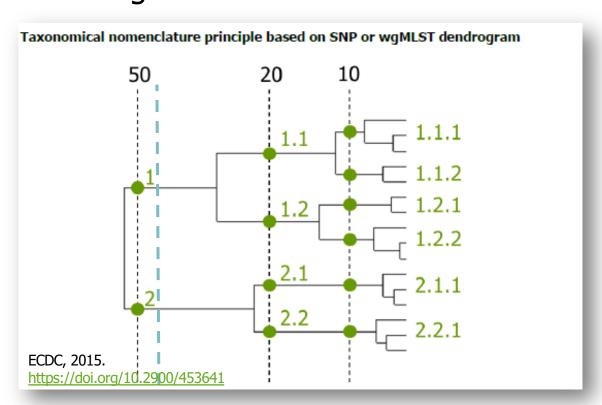
S - singleton



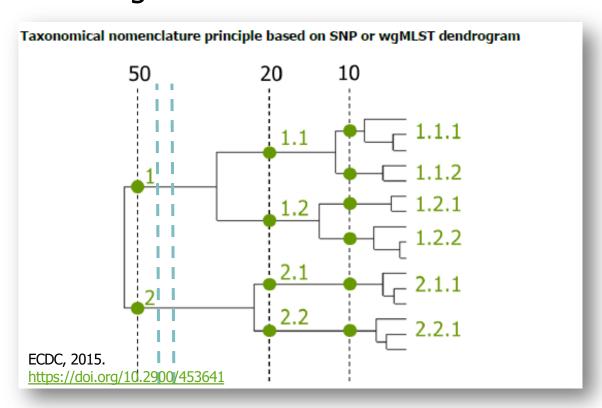
# How can I determine the nomenclature code levels?



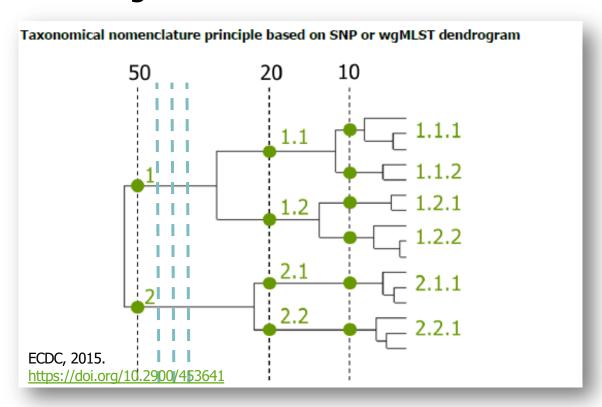




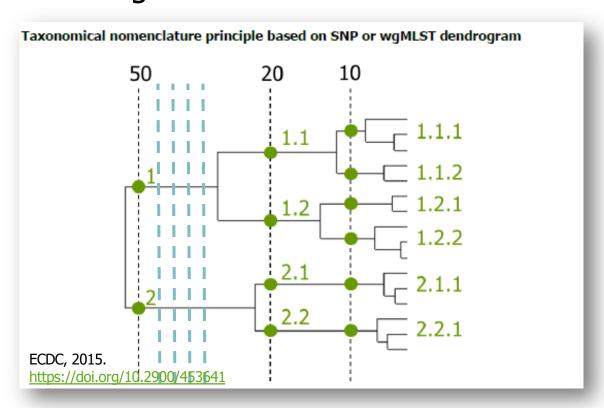






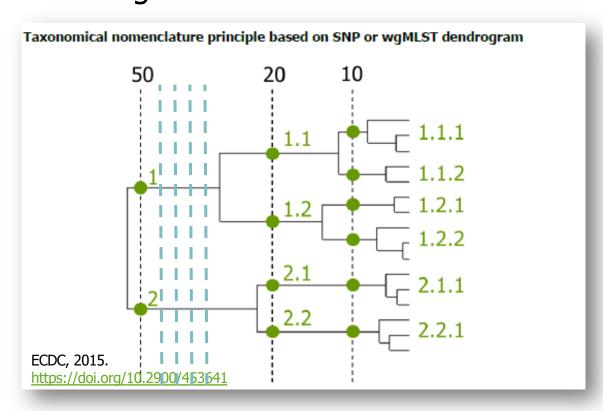








Stability regions are ranges of subsequent distance thresholds providing similar clustering results.



Candidate regions to define cutoffs for pathogen-specific nomenclature design, as they are less likely to be "disrupted" by the introduction of new samples

### Identification of candidate thresholds for nomenclature design



Neighborhood Adjusted Wallace coefficient (nAWC) – evaluates the agreement between two clustering methods (e.g. subsequent distance thresholds)

### Identification of candidate thresholds for nomenclature design

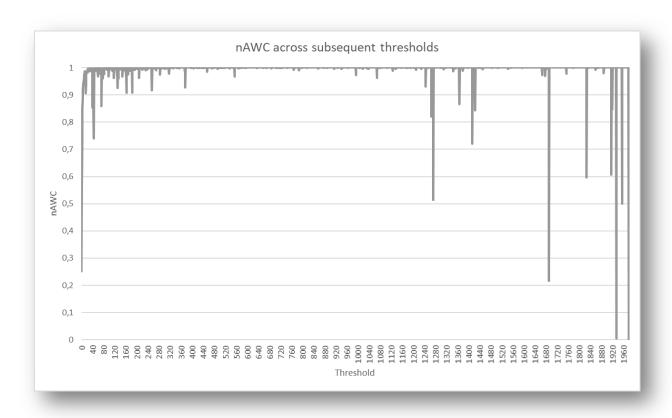


Neighborhood Adjusted Wallace coefficient (nAWC) – evaluates the agreement between two clustering methods (e.g. subsequent distance thresholds)

comparing\_partitions\_v2.py

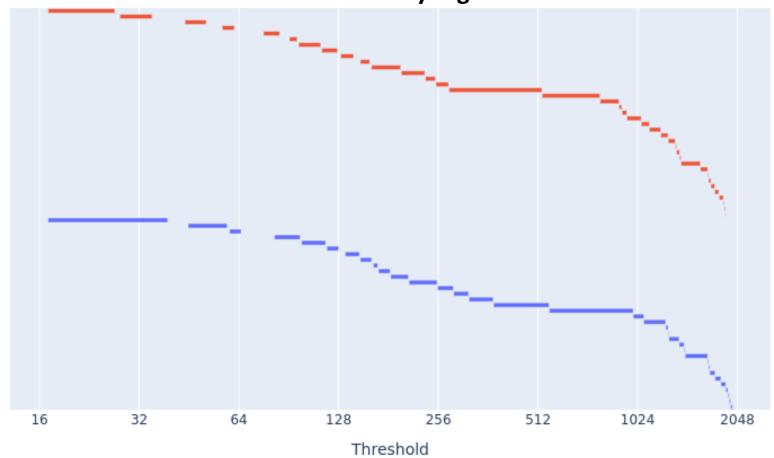


Determines the nAWC at consecutive thresholds  $("n + 1" \rightarrow "n")$ 





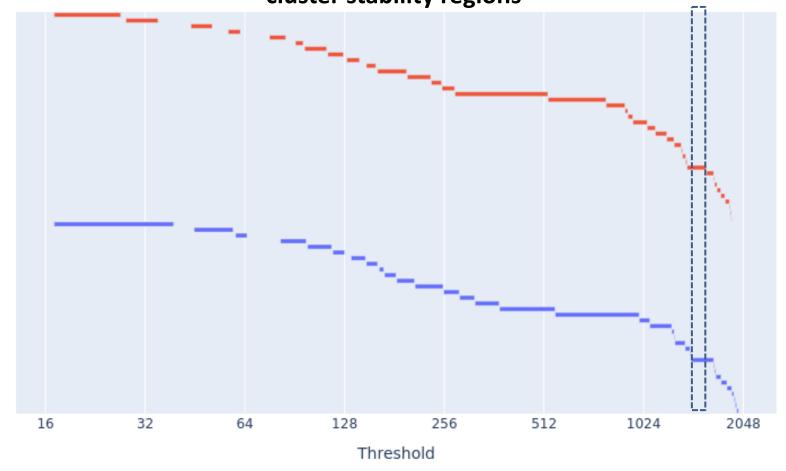
Threshold ranges in which clustering remains similar\*
- cluster stability regions -



<sup>\*</sup>Blocks identified with comparing\_partitions\_v2.py for two pipelines running the same dataset (visualization: EvalTree module of ACDTree toolkit)



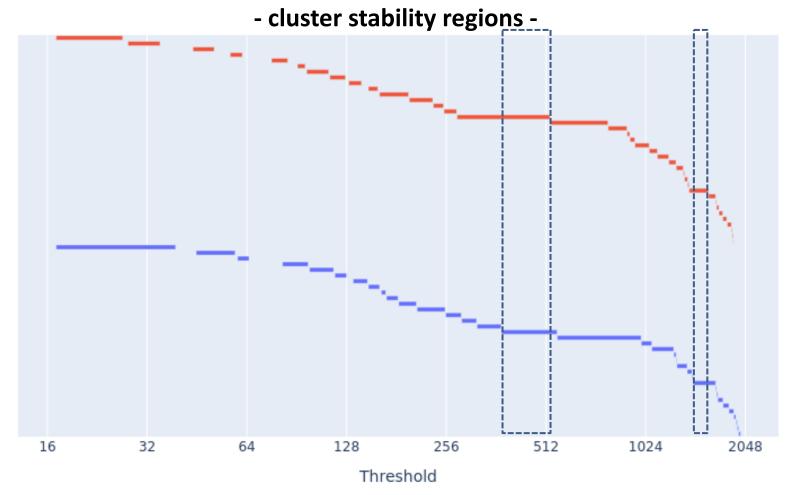
Threshold ranges in which clustering remains similar\*
- cluster stability regions -



<sup>\*</sup>Blocks identified with comparing\_partitions\_v2.py for two pipelines running the same dataset (visualization: EvalTree module of ACDTree toolkit)



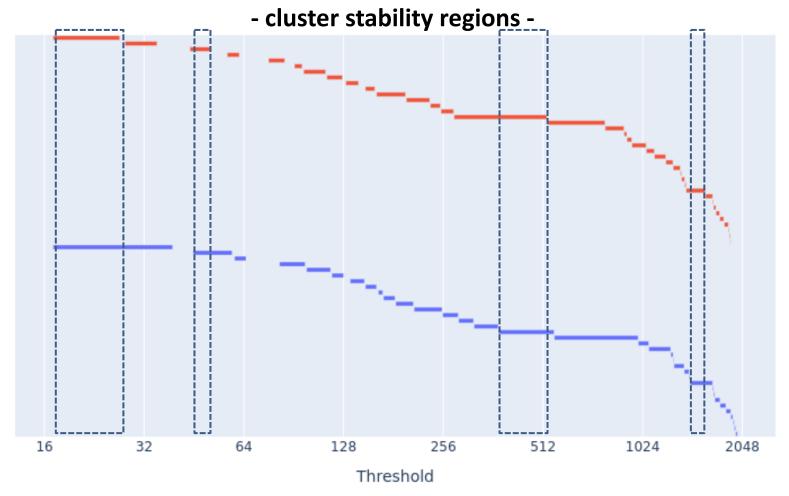
Threshold ranges in which clustering remains similar\*



<sup>\*</sup>Blocks identified with comparing\_partitions\_v2.py for two pipelines running the same dataset (visualization: EvalTree module of ACDTree toolkit)



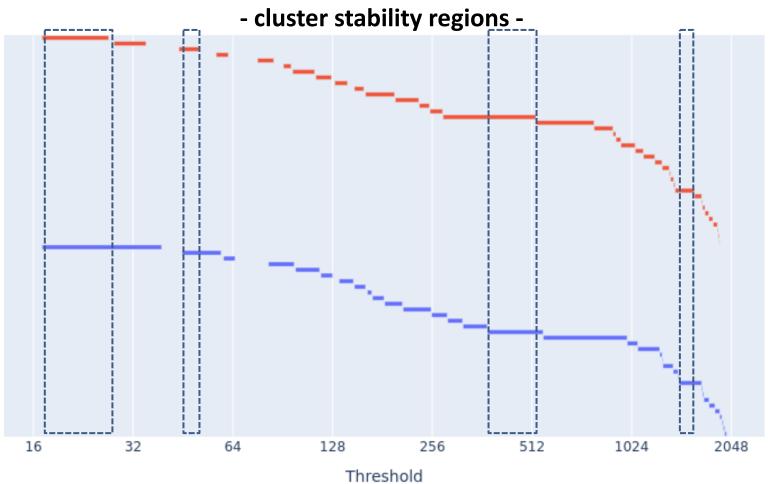
Threshold ranges in which clustering remains similar\*



<sup>\*</sup>Blocks identified with comparing\_partitions\_v2.py for two pipelines running the same dataset (visualization: EvalTree module of ACDTree toolkit)



Threshold ranges in which clustering remains similar\*



--nomenclature-code-levels 1500,500,50,20

<sup>\*</sup>Blocks identified with comparing\_partitions\_v2.py for two pipelines running the same dataset (visualization: EvalTree module of ACDTree toolkit)

### Cluster stability regions across pipelines and WGS backward compatibility with other typing solutions

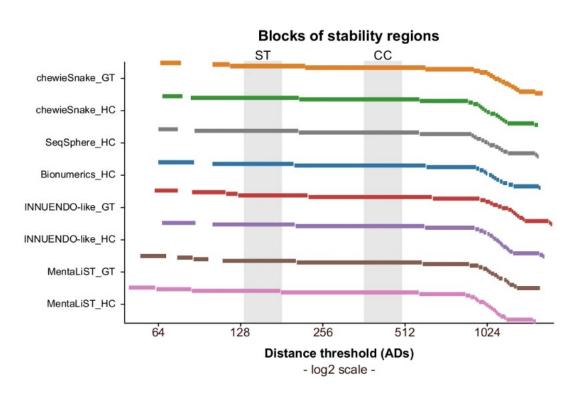


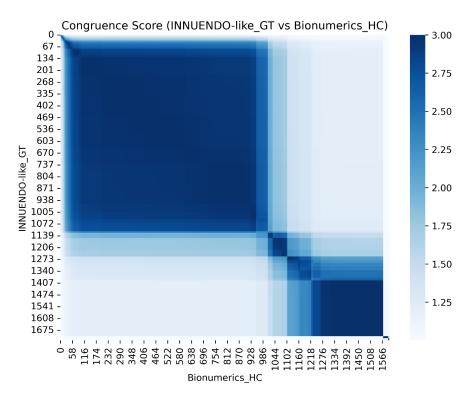
• As they reflect the population structure, some cluster stability regions may correspond to some genetic-based traditional typing -> this can also be a criteria to use it in the nomenclature code

### Cluster stability regions across pipelines and WGS backward compatibility with other typing solutions



• As they reflect the population structure, some cluster stability regions may correspond to some genetic-based traditional typing -> this can also be a criteria to use it in the nomenclature code



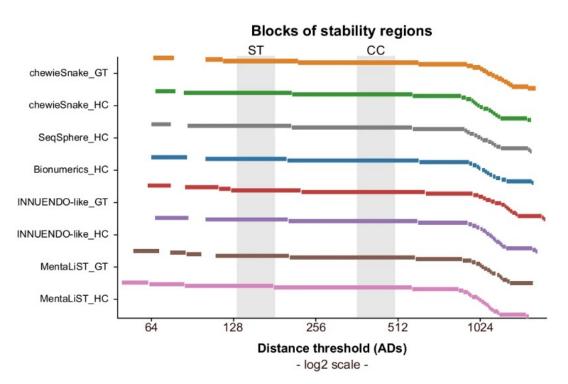


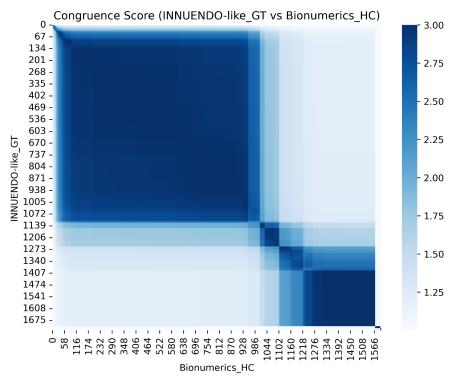
E.g. large-scale multicountry inter-pipeline cluster congruence (Food- and water-borne bacterial pathogens)

### Cluster stability regions across pipelines and WGS backward compatibility with other typing solutions



• As they reflect the population structure, some cluster stability regions may correspond to some genetic-based traditional typing -> this can also be a criteria to use it in the nomenclature code





Given its input flexibility and by providing clustering information at all possible threshold levels, ReporTree facilitates large-scale congruence analyses!

#### Exercise – genomic surveillance of Salmonella enterica



- (Pre-webinar: install ReporTree)
- Download the provided material (reportree\_exercise.zip) and paste it in ReporTree folder
- unzip reportree exercise.zip
- Check its contents (ls reportree\_exercise/):
  - o run last week/ folder with the results of a previous ReporTree run
  - o run\_this\_week/ folder with the input files necessary for this ReporTree run:

Se\_beone\_alleles.tsv - wgMLST allele matrix for all samples in the dataset (including

samples added since "last week")

Se beone metadata.tsv — metadata table for all samples in the dataset

cgMLST\_Salmonella\_enterica\_chewie-NS.txt - List of cgMLST loci

#### Exercise – genomic surveillance of Salmonella enterica



#### **Genomic surveillance command line:**

```
python /PATH/TO/reportree.py -a Se beone alleles.tsv \
                             -m Se beone metadata.tsv \
                             -out Se beone GT \
                             -1 cgMLST Salmonella enterica chewie-NS.txt
                             --loci-called 0.95 \
                             -thr 1,2,5,10,25,100,500,1000 \
                             --analysis grapetree \
                             --nomenclature-file ../run last week/Se beone last week GT partitions.tsv \
                             --nomenclature-code 500,100,25,10 \
                             --columns_summary_report MLST_ST, Predicted_serotype, Country, n_Country, Source,
n Source, first seq date, last seq date, timespan days
```

#### Exercise – genomic surveillance of Salmonella enterica



#### **Genomic surveillance command line:**

```
python /PATH/TO/reportree.py -a Se beone alleles.tsv \
                             -m Se beone metadata.tsv \
                              -out Se beone GT \
                              -1 cgMLST Salmonella enterica chewie-NS.txt
                              --loci-called 0.95 \
                              -thr 1,2,5,10,25,100,500,1000 \
                              --analysis grapetree \
                              --nomenclature-file ../run last week/Se beone last week GT partitions.tsv \
                              --nomenclature-code 500,100,25,10 \
                              --columns_summary_report MLST_ST, Predicted_serotype, Country, n_Country, Source,
n Source, first seq date, last seq date, timespan days
+ Zoom-in analysis: --sample of interest ERR10442916 --zoom-cluster-of-interest 5 --site-inclusion 1.0 --unzip
+ Finding cluster stability regions: --partitions2report stability regions -thr-1,2,5,10,25
```





List of learning points in this session:

ReporTree is a <u>flexible solution that automatically identifies genetic clusters</u> at any (or all)
distance thresholds and generates surveillance-oriented reports



- ReporTree is a <u>flexible solution that automatically identifies genetic clusters</u> at any (or all)
   distance thresholds and generates surveillance-oriented reports
- It can be smoothly implemented in **routine surveillance**, providing a <u>cluster nomenclature</u> <u>system</u> that facilitates monitoring main circulating lineages and ongoing outbreaks



- ReporTree is a <u>flexible solution that automatically identifies genetic clusters</u> at any (or all) distance thresholds and generates surveillance-oriented reports
- It can be smoothly implemented in **routine surveillance**, providing a <u>cluster nomenclature</u> <u>system</u> that facilitates monitoring main circulating lineages and ongoing outbreaks
- When dealing with an wgMLST schema, ReporTree automatically performs a a <u>dynamic</u> wgMLST approach to zoom-in in any cgMLST cluster of interest



- ReporTree is a <u>flexible solution that automatically identifies genetic clusters</u> at any (or all)
  distance thresholds and generates surveillance-oriented reports
- It can be smoothly implemented in **routine surveillance**, providing a <u>cluster nomenclature</u> <u>system</u> that facilitates monitoring main circulating lineages and ongoing outbreaks
- When dealing with an wgMLST schema, ReporTree automatically performs a a <u>dynamic</u> wgMLST approach to zoom-in in any cgMLST cluster of interest
- As a complement, it can identify <u>clustering stability regions</u>, a useful approach for pathogen specific <u>nomenclature design</u>



List of learning points in this session:

- ReporTree is a <u>flexible solution that automatically identifies genetic clusters</u> at any (or all) distance thresholds and generates surveillance-oriented reports
- It can be smoothly implemented in **routine surveillance**, providing a <u>cluster nomenclature</u> <u>system</u> that facilitates monitoring main circulating lineages and ongoing outbreaks
- When dealing with an wgMLST schema, ReporTree automatically performs a a <u>dynamic</u> wgMLST approach to zoom-in in any cgMLST cluster of interest
- As a complement, it can identify <u>clustering stability regions</u>, a useful approach for pathogen specific <u>nomenclature design</u>

ReporTree contributes to a sustainable and efficient public health genomics-informed pathogen surveillance

### **Further reading**



Stay tuned!
Updates coming soon!!!

#### ReporTree GitHub:

https://github.com/insapathogenomics/ReporTree

#### Clustering stability regions:

https://github.com/insapathogenomics/ComparingPartitions

https://github.com/insapathogenomics/ACDTree





#### References



- Dallman et al. (2018) SnapperDB: a database solution for routine sequencing analysis of bacterial isolates. *Bioinformatics*. 34(17):3028-3029. doi: 10.1093/bioinformatics/bty212
- ECDC (2025) One Health (<a href="https://www.ecdc.europa.eu/en/one-health">https://www.ecdc.europa.eu/en/one-health</a>)
- Francisco et al. (2012) PHYLOViZ: phylogenetic inference and data visualization for sequence based typing methods. BMC Bioinformatics. 13:87. doi: 10.1186/1471-2105-13-87
- Mixão et al. (2023) ReporTree: a surveillance-oriented tool to strengthen the linkage between pathogen genetic clusters and epidemiological data. *Genome Med.* 15(1):43. doi: 10.1186/s13073-023-01196-1
- Mixão et al. (2025) Multi-country and intersectoral assessment of cluster congruence between pipelines for genomics surveillance of foodborne pathogens. *Nat Commun.* 16(1):3961. doi: 10.1038/s41467-025-59246-8
- Severiano et al. (2011) Adjusted Wallace coefficient as a measure of congruence between typing methods. J Clin Microbiol. 49(11):3997-4000. doi: 10.1128/JCM.00624-11
- WHO (2022) WHO global genomic surveillance strategy for pathogens with pandemic and epidemic potential 2022-2032 (<a href="https://www.who.int/initiatives/genomic-surveillance-strategy">https://www.who.int/initiatives/genomic-surveillance-strategy</a>)
- Zhou et al. (2018) GrapeTree: visualization of core genomic relationships among 100,000 bacterial pathogens. Genome Res. 28(9):1395-1404. doi: 10.1101/gr.232397.117



### Genomics and Bioinformatics Unit of the National Institute of Health Doutor Ricardo Jorge (INSA)

Daniel Sobral
Miguel Pinto
João Paulo Gomes
<u>Vítor Borges</u>

### Instituto\_Nacional de Saúde Doutor Ricardo Jorge











veronica.mixao@insa.min-saude.pt

### **Acknowledgements**

The creation of this training material was commissioned by ECDC to National Institute of Health Doutor Ricardo Jorge (INSA) with the direct involvement of Verónica Mixão.

ReporTree development was supported by co-funding from the European Union's Horizon 2020 Research and Innovation program under grant agreement No 773830: One Health European Joint Programme.

ReporTree expansion and continuous maintenance is supported by national funds through FCT - Foundation for Science and Technology, I.P., in the frame of Individual CEEC 2022.00851.CEECIND/CP1748/CT0001, by the European Union project "Sustainable use and integration of enhanced infrastructure into routine genome-based surveillance and outbreak investigation activities in Portugal" - GENEO [101113460] on behalf of the EU4H programme [EU4H-2022-DGA-MS-IBA-01-02], and by the DURABLE "Research Network against Epidemics" project. DURABLE is co-funded by The European Commission Union under the EU4Health Programme (EU4H) [101102733]. However, views and opinions expressed are those of the authors only and do not necessarily reflect those of the European Union or the European Health and Digital Executive Agency. Neither the European Union nor the granting authority can be held responsible.