



**Outbreaks, typing & AMR/Day 7**

# **Serotyping and virulence typing of *Salmonella* and *E. coli***

Section of Foodborne infections

Dep. of Bacteria, Parasites and Fungi, Statens Serum Institut

Katrine Grimstrup Joensen and Pernille Gymoese, April 2024

# Objectives

Specific objectives of this session:

1. Learn about *in silico* serotyping of *Salmonella* and *E. coli* isolates
2. Learn about *in silico* virulence profiling of *E. coli*
3. Learn how to apply SerotypeFinder, SeqSero and VirulenceFinder
4. Interpret typing results

Related to the course objectives:

- A. Process sequencing data (from wgs data to results)
- B. Extract relevant information from processed data
- C. Perform basic analysis supporting epidemiological investigations, including interaction with public databases

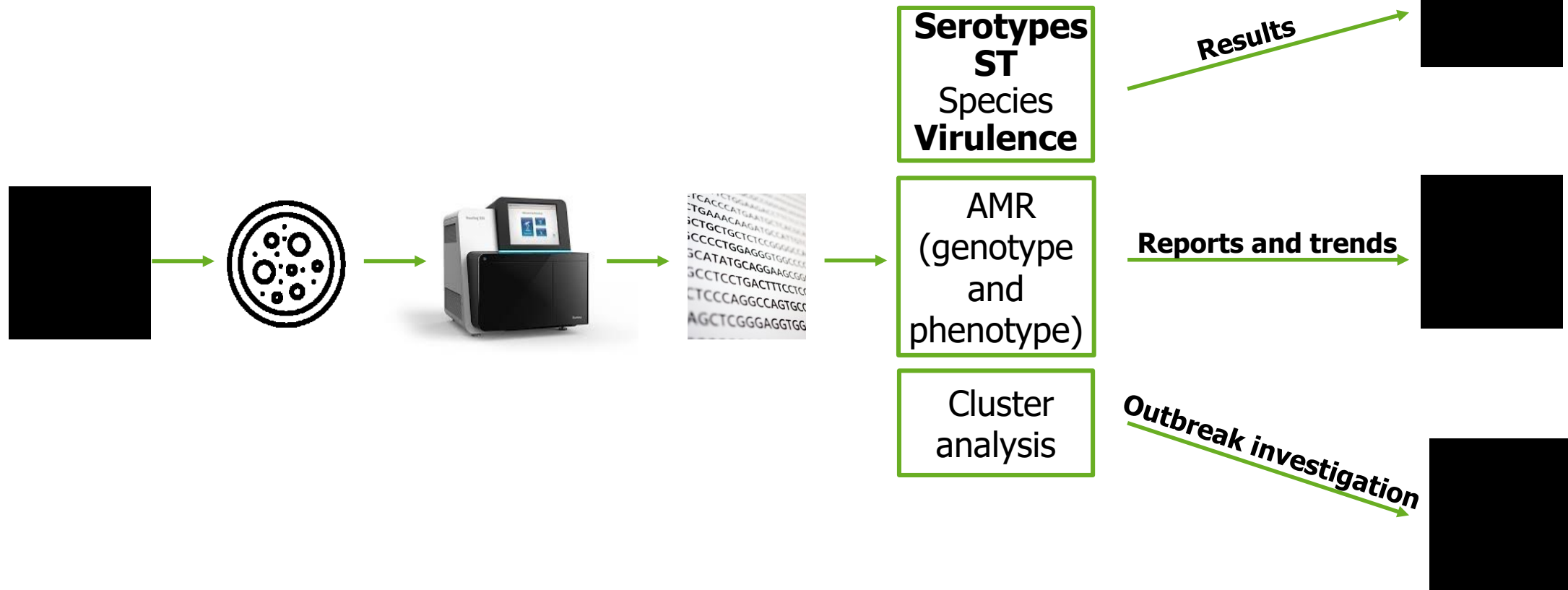
# Outline

This session consists of the following elements

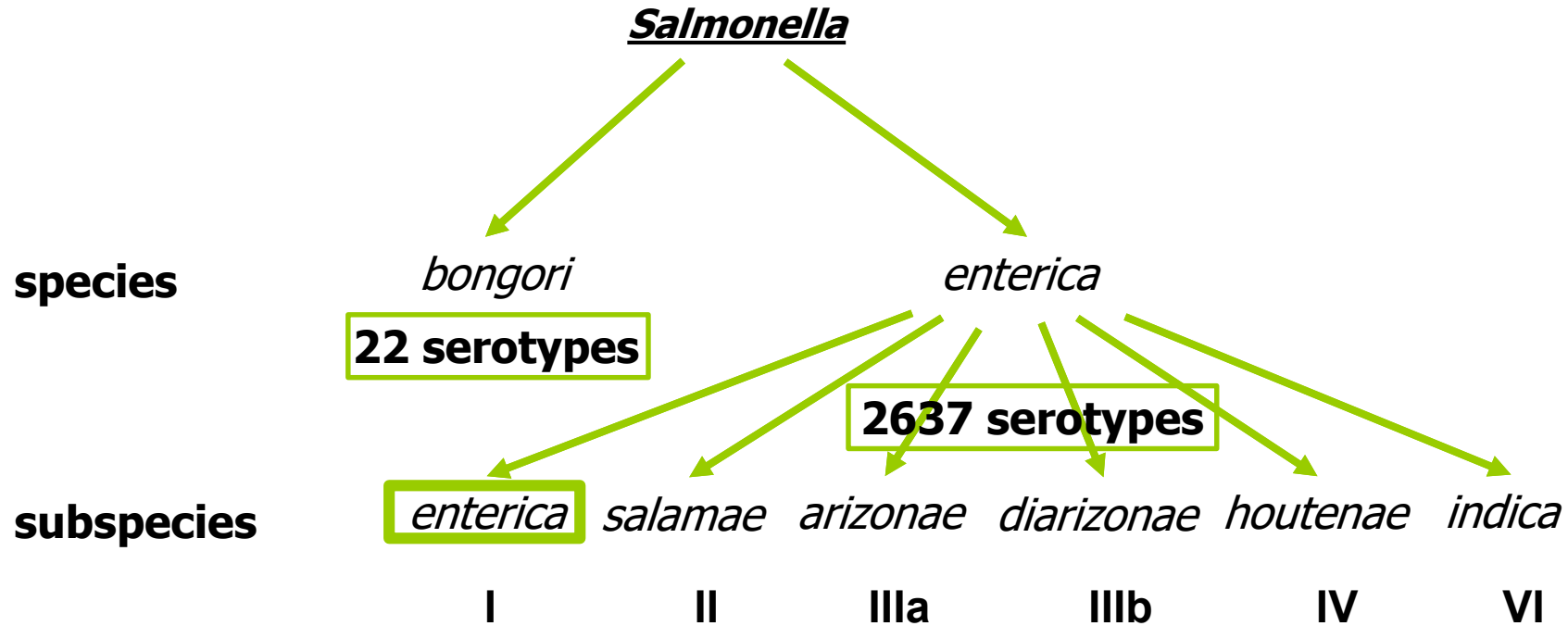
1. Introduction to *in silico* serotyping of *Salmonella*
2. SeqSero
3. Introduction to *in silico* serotyping and virulence profiling of *E. coli*
4. SerotypeFinder
5. VirulenceFinder
6. Issues regarding *in silico* typing
7. Introduction to exercise
8. Group exercise integrating *in silico* tools for serotyping of *E. coli* and *Salmonella* and virulence profiling of *E. coli*
9. Group discussion of results

# Surveillance at SSI

Surveillance flow from hospital to SSI and back



# Serotyping of *Salmonella*



**White-Kauffmann-Le Minor scheme**  
O and H antigens

Host, disease, place of isolation

- Based on O and H antigens
- **O:H phase 1:H phase 2**
- Subspecies *enterica* is named after host, disease or geographical place of isolation
- Other subspecies after antigen formula

## Examples:

- **Serotype:** *Salmonella enterica* subspecies *enterica* serotype London (*S. London*)
- **Antigen profile:** 3,10:l,v:1,6
- **Serotype:** *Salmonella enterica* subspecies *salamae* serotype II 58:a:z<sub>6</sub>

# Sequence typing – new nomenclature

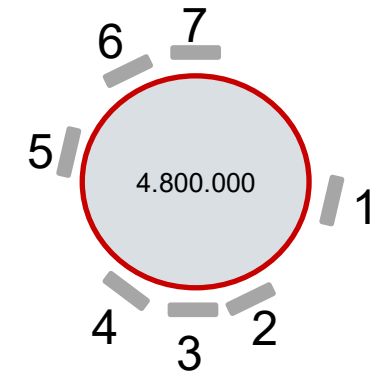
- **Multi Locus Sequence Typing - MLST**

- Seven conserved genes
- Genetic relatedness

- **eBG – eBurstGroup**

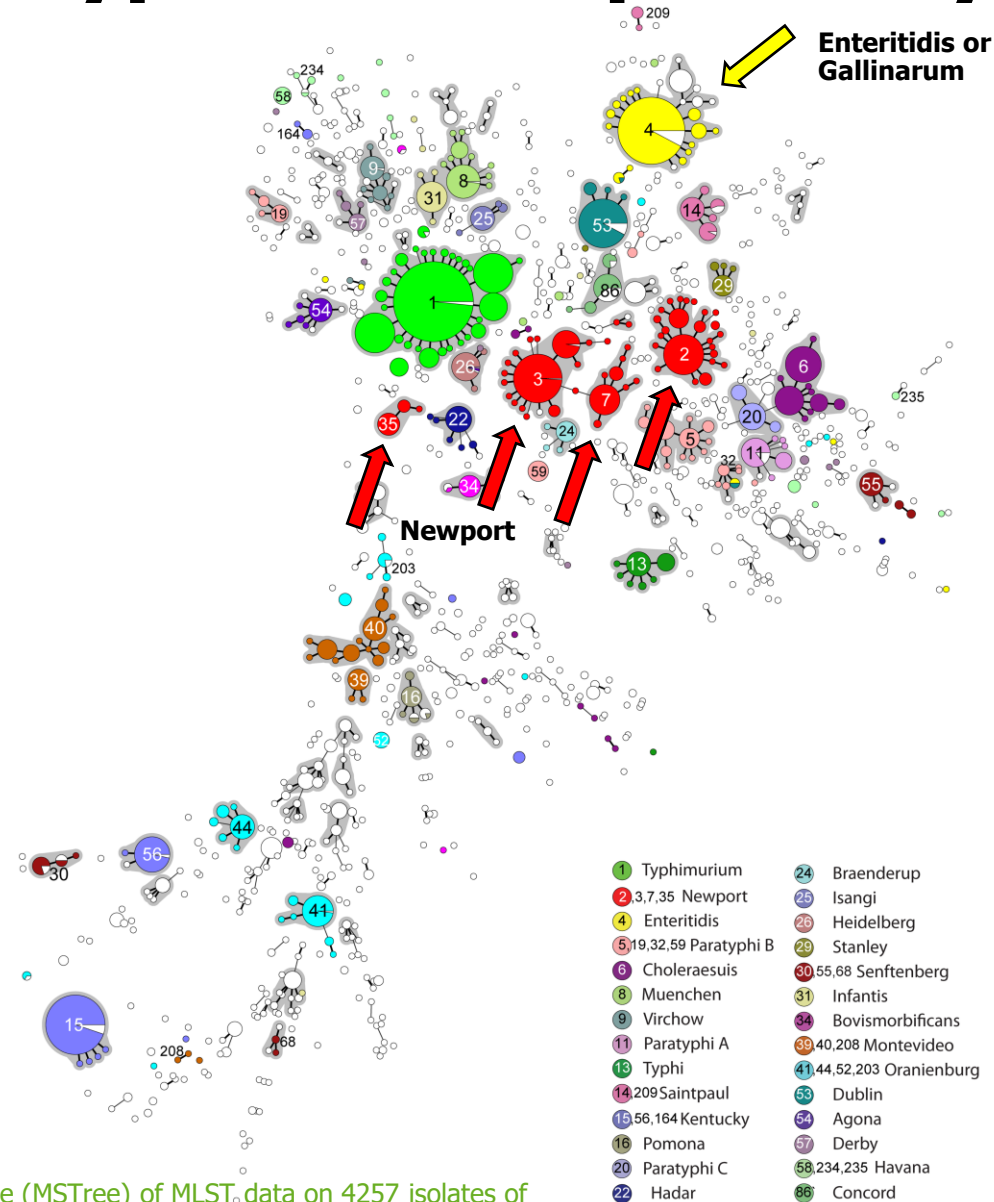
- Closely related ST's
- Correlates to serotypes

MLST - 7 genes



aroC	dnaN	hemD	hisD	purE	sucA	thrA		ST	eBG	Serotype
17	18	22	17	5	21	19		32	31	Infantis
17	341	22	17	5	21	19		1032	31	Infantis

# Serotypes and sequence types



## Serotypes

- Pros
  - Well known nomenclature
  - Can provide information on properties like disease and host specificity
- Cons
  - Not necessarily genetic related

## ST and eBG

- Pros
  - Information on genetic relatedness
- Cons
  - No information on properties

**Replacement for serotyping? – or a combination of both**

# Interpretation of *Salmonella* serotype results

- **Same antigen formula - different subspecies**
- **Biotypes/variants** – same antigen formula – different phenotypes
  - Additional biochemical test
- **Paratyphi B and Paratyphi B var. Java**
  - Differentiated by fermentation of d-tartrate
- **Enteritidis/Gallinarum**
  - Genotypic serotype looks like Enteritidis – non-expressed H phase
  - Phenotypic serotype looks like non-motile Enteritidis
  - Additional biochemical test
  - Gene *sdf* missing in Gallinarum and present in common Enteritidis – but is not exclusive....
- See <http://denglab.info/SeqSero2/supple> or [https://www.pasteur.fr/sites/default/files/veng\\_0.pdf](https://www.pasteur.fr/sites/default/files/veng_0.pdf) for further details



# SeqSero2

- **Input**
  - Raw reads
  - Assemblies
- **Output**
  - Detected O and H antigenes
  - Predicted antigen formula
  - Predicted serotype
  - Predicted species and subspecies
  - Potential inter-serotype contamination
  - Notes

# SeqSero2 vs SeqSero1

- **Faster – k-mer workflow**
- **Identify species and subspecies**
- **Additional genetic markers**
  - Differentiate Paratyphi B and variant Java (D-tartrate fermentation)
  - +/- *sdf* gene in Enteritidis – indication of variant Gallinarum (*sdf* negative)
  - O5 negative mutation (-O5) – variant of Typhimurium
  - Differentiation of O22 and O23 in the O13 group
- **Inter-serotype contamination**
  - Can determine more than one O antigen and more than 2 H antigens – indication of contamination

# Alternatives

- **SISTR**

- Serotype prediction from MLST and extended MLST schemes (rMLST and cgMLST)
- <https://doi.org/10.1371/journal.pone.0147101>

- **Enterobase**

- Seqsero2, SISTR and predicted type after ST/eBG
- <https://enterobase.warwick.ac.uk/>



The screenshot shows the Enterobase website interface. At the top, it says "Enterobase A powerful, user-friendly online resource for analysing and visualising genomic variation within enteric bacteria". Below this, there are navigation links for "Log In", "Register", "Help", and "v1.2.0". The main content area displays "Enterobase currently contains 1,032,631 bacterial strains". There are six panels, each representing a different bacterial genus with a background image of the bacteria and a text box containing statistics and scheme information:

- Salmonella**: Total strains: 453,770; Assemblies in progress: 18; Schemes: rMLST, Achtman 7 Gene MLST, cgMLST V2 + HierCC V1, wgMLST
- Escherichia/Shigella**: Total strains: 294,380; Assemblies in progress: 15; Schemes: rMLST, wgMLST, Achtman 7 Gene MLST, cgMLST V1 + HierCC V1
- Mycobacterium**: Total strains: 113,995; Assemblies in progress: 3
- Streptococcus**: Total strains: 99,600; Assemblies in progress: 8; Schemes: cgMLST v1
- Clostridioides**: Total strains: 30,388; Assemblies in progress: 1; Schemes: wgMLST, rMLST, cgMLST V1 + HierCC V1
- Vibrio**: Total strains: 17,377; Assemblies in progress: 0; Schemes: rMLST, Vibrio cgMLST + HierCC V1

# Serotyping of *E. coli*

- Based on O and H antigens
- Typical single O and H antigens
- Named according to antigen formula

## Example:

- O157:H7

# *E. coli* pathotypes

- **Diarrheagenic *E. coli* (DEC)**
  - **STEC - Shiga toxin-producing**
  - EPEC - Enteropathogenic
  - ETEC - Enterotoxigenic
  - EIEC - Enteroinvasive
  - EAEC - Enteroaggregative
- **Extraintestinal *E. coli* (ExPEC)**
  - UPEC - Uropathogenic

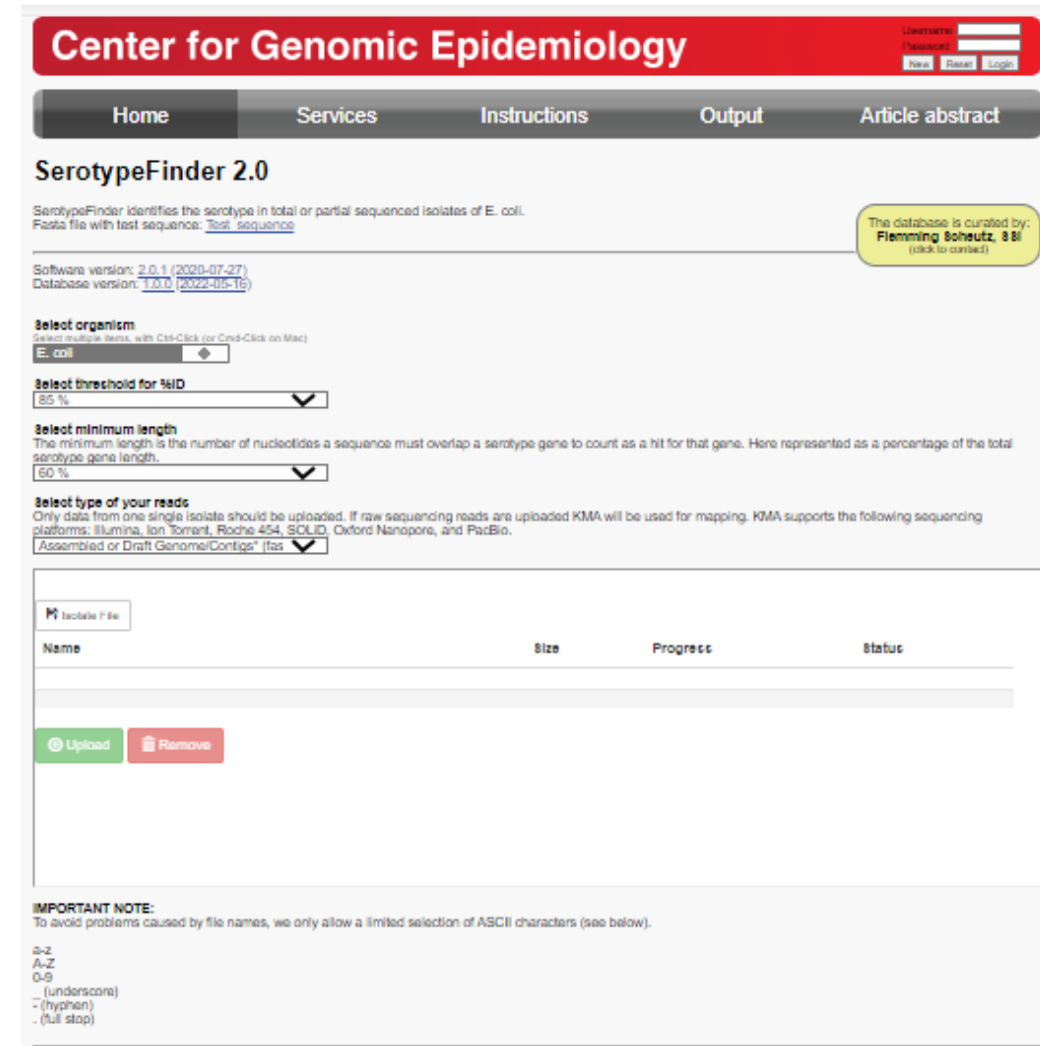
Characterized by  
virulence genes

# *E. coli* typing

- **Previous STEC typing at SSI:**
  - Serotyping
  - Resistance profiling
  - Test for hemolysin,  $\beta$ -glucuronidase, indole production
  - PCR and Hybridization to detect virulence genes (*eae*, *vtx1*, *vtx2*...)
  - PFGE for suspected outbreak isolates
- **Replaced by WGS analysis**

# SerotypeFinder & VirulenceFinder

- **Developed as online CGE tools**
  - Select database
  - Select ID/Coverage (default 85%, 60%)
  - Select data type
  - Upload data
- **Blast-based, KMA for raw data**
- **Also available command-line**
- **At SSI - integrated databases (not the whole tool) in our pipeline (with mapping)**



**Center for Genomic Epidemiology**

Home Services Instructions Output Article abstract

### SerotypeFinder 2.0

SerotypeFinder identifies the serotype in total or partial sequenced isolates of E. coli.  
Fasta file with first sequence: [Test sequence](#)

The database is curated by: **Flemming Boetzel, 881** (click to contact)

Software version: [2.0.1 \(2020-07-27\)](#)  
Database version: [1.0.0 \(2022-05-18\)](#)

**Select organism**  
Select multiple items with Ctrl-Click (or Cmd-Click on Mac)  
E. coli

**Select threshold for %ID**  
85 %

**Select minimum length**  
The minimum length is the number of nucleotides a sequence must overlap a serotype gene to count as a hit for that gene. Here represented as a percentage of the total serotype gene length.  
60 %

**Select type of your reads**  
Only data from one single isolate should be uploaded. If raw sequencing reads are uploaded KMA will be used for mapping. KMA supports the following sequencing platforms: Illumina, Ion Torrent, Roche 454, SOLiD, Oxford Nanopore, and PacBio.  
Assembled or Draft Genome/Contigs (fas)

Upload File

Name	Size	Progress	Status
------	------	----------	--------

Upload Remove

**IMPORTANT NOTE:**  
To avoid problems caused by file names, we only allow a limited selection of ASCII characters (see below).  
a-z  
A-Z  
0-9  
\_ (underscore)  
- (hyphen)  
. (full stop)

# SerotypeFinder: *in silico* O and H-typing

- **187 known *E. coli* O-antigens (O1-O187)**
  - More added now
  - ***wzx*** and ***wzy*** (O-antigen flippase and polymerase)
  - ***wzm*** and ***wzt*** (ABC transporter)
- **All 53 known H-types**
  - *fliC, flkA, flnA, flmA, fliA*



# SerotypeFinder results

- **Input**

- Raw reads
- Assemblies

- **Output**

- O - combined wzx/wzy
- H - flagellin genes
- Sequence identity, coverage and positions
- Reference sequence
- Examine alignments
- Command line output more messy and needs to be sorted

## SerotypeFinder-2.0 Server - Results

Database(s): *O\_type,H\_type*

Database for O type genes						
Gene	Serotype	Identity	Template / HSP length	Contig	Position in contig	Accession number
wzy	O104	99.64	1113 / 1113	NODE_122_length_26982_cov_17.237417	16386..17498	AF361371
wzx	O104	100	1278 / 1278	NODE_122_length_26982_cov_17.237417	18727..20004	KB021482

Database for H type genes						
Gene	Serotype	Identity	Template / HSP length	Contig	Position in contig	Accession number
flhC	H4	100	1050 / 1050	NODE_24_length_15621_cov_15.220664	9868..10917	AJ605764

extended output

Results as text Results tsv Hits in genome seqs Serotype gene sequences

Selected %ID threshold: 85 %

Selected minimum length: 60 %

Input Files: *test.txt*

# VirulenceFinder

- Initially 76 virulence genes – and growing....

## Center for Genomic Epidemiology

Username   
Password

Home Services Instructions Output Article abstract

### VirulenceFinder 1.4

Course in Whole Genome Sequencing and Analysis for Clinical Microbiologists, click [here](#) for more information  
View the [version history](#) of this server.

The database is curated by:  
Flemming Scheutz, SSI  
(click to contact)

Select species  
E. coli  
Enterococcus  
S. aureus

Select threshold for %ID  
98 %

Select type of your reads  
Assembled Genome/Contigs\*

Isolate File

Name	Size	Progress	Status

\* Please note also that "Assembled Genomes/Contigs" should be selected, if you have already assembled your short sequencing reads into one continuous genome or into several contigs. "Assembled Genomes/Contigs" is defined as one or several contigs in one FASTA file (one entry per contig). It is indifferent which type of short sequence reads were used to produce the genome/contigs.



## Real-Time Whole-Genome Sequencing for Routine Typing, Surveillance, and Outbreak Detection of Verotoxigenic *Escherichia coli*

Katrine Grimstrup Joensen,<sup>a,b</sup> Flemming Scheutz,<sup>b</sup> Ole Lund,<sup>c</sup> Henrik Hasman,<sup>a</sup> Rolf S. Kaas,<sup>a,c</sup> Eva M. Nielsen,<sup>d</sup> Frank M. Aarestrup<sup>a</sup>  
National Food Institute, Division for Epidemiology and Microbial Genomics, Technical University of Denmark, Kongens Lyngby, Denmark<sup>a</sup>; Department of Microbiology and Infection Control, Statens Serum Institut, Copenhagen, Denmark<sup>b</sup>; Center for Biological Sequence Analysis, Department of System Biology, Technical University of Denmark, Kongens Lyngby, Denmark<sup>c</sup>

Fast and accurate identification and typing of pathogens are essential for effective surveillance and outbreak detection. The current routine procedure is based on a variety of techniques, making the procedure laborious, time-consuming, and expensive. With whole-genome sequencing (WGS) becoming cheaper, it has huge potential in both diagnostics and routine surveillance. The aim of this study was to perform a real-time evaluation of WGS for routine typing and surveillance of verocytotoxin-producing *Escherichia coli* (VTEC). In Denmark, the Statens Serum Institut (SSI) routinely receives all suspected VTEC isolates. During a 7-week period in the fall of 2012, all incoming isolates were concurrently subjected to WGS using IonTorrent PGM. Real-time bioinformatics analysis was performed using web-tools ([www.genomicepidemiology.org](http://www.genomicepidemiology.org)) for species determination, multilocus sequence type (MLST) typing, and determination of phylogenetic relationship, and a specific VirulenceFinder for detection of *E. coli* virulence genes was developed as part of this study. In total, 46 suspected VTEC isolates were characterized in parallel during the study. VirulenceFinder proved successful in detecting virulence genes included in routine typing, explicitly verocytotoxin 1 (*vtx1*), verocytotoxin 2 (*vtx2*), and intimin (*eae*), and also detected additional virulence genes. VirulenceFinder is also a robust method for assigning verocytotoxin (*vtx*) subtypes. A real-time clustering of isolates in agreement with the epidemiology was established from WGS, enabling discrimination between sporadic and outbreak isolates. Overall, WGS typing produced results faster and at a lower cost than the current routine. Therefore, WGS typing is a superior alternative to conventional typing strategies. This approach may also be applied to typing and surveillance of other pathogens.

# VirulenceFinder

- **Detection of pathotype-specific genes**
- **STEC**
  - *eae* - 45 variants
  - *stx1* - 14A,18B variants
  - *stx2* - 114A,43B variants
  - *stx* subtypes
    - Important for HUS associated types

Gene	Description	Variants in database
<i>astA</i>	Heat-stable enterotoxin 1	11
<i>bfpA</i>	Major subunit of bundle-forming pili	5
<i>cba</i>	Colicin B	15
<i>ccl</i>	Cloacin	4
<i>cdtB</i>	Cytolethal distending toxin B	14
<i>celB</i>	Endonuclease colicin E2	10
<i>cfa_c</i>	Colonisation factor antigen I	4
<i>cif</i>	Type III secreted effector	4
<i>cma</i>	Colicin M	19
<i>cnf1</i>	Cytotoxic necrotizing factor	7
<i>cofA</i>	Longus type IV pilus subunit	1
<i>eae</i>	Intimin	45
<i>eatA</i>	SPATE	3
<i>efa1</i>	EHEC factor for adherence	11
<i>ehxA</i>	Enterohaemolysin	12
<i>epeA</i>	SPATE	1
<i>espA</i>	Type III secretions system	23
<i>espB</i>	Secreted protein B	14
<i>espC</i>	SPATE	3
<i>espF</i>	Type III secretion system	13
<i>espI</i>	SPATE	2
<i>espJ</i>	Prophage-encoded t3SS effector	2
<i>espP</i>	Putative exoprotein precursor	4
<i>etpD</i>	Type II secretion protein	3
<i>f17-A</i>	Subunit A of F17 fimbrial protein	7
<i>f17-G</i>	Adhesin subunit of F17 fimbriae	9
<i>fanA</i>	Involved in biogenesis of K99/F5 fimbriae	1
<i>fasA</i>	Fimbriae 987P/F6 subunit	1
<i>fedA</i>	Fimbrial protein F107 subunit A	3
<i>fedF</i>	Fimbrial adhesin AC precursor	6
<i>fim41a</i>	Mature Fim41a/F41 protein	2
<i>gad</i>	Glutamate decarboxylase	70
<i>hlyE</i>	Avian <i>E.coli</i> haemolysin	1
<i>iha</i>	Adherence protein	19
<i>ipaD</i>	Invasion protein <i>Shigella flexneri</i>	9
<i>ipaH9.8</i>	Invasion plasmid antigen	8
<i>ireA</i>	Siderophore receptor	4
<i>iroN</i>	Enterobactin siderophore receptor	13

# VirulenceFinder results

- **Input**

- Raw reads
- Assemblies

- **Output**

- Sequence identity, coverage and positions
- Reference sequence
- Examine alignments
- Sort in relevance of genes

## VirulenceFinder-2.0 Server - Results

5315833.cgebase.food.dtu.dk

Organism(s): *Escherichia coli*

Shiga-toxin genes						
Virulence factor	Identity	Query / Template length	Contig	Position in contig	Protein function	Accession number
stx2	100	1241 / 1241	NODE_590_length_1768_cov_60.913460	438..1678	O157 SF-3573-98, variant a	<a href="#">AB030484</a>

Virulence genes for Escherichia coli						
Virulence factor	Identity	Query / Template length	Contig	Position in contig	Protein function	Accession number
ORF3	100	1029 / 1029	NODE_27_length_2659_cov_706.860840	348..1376	Isoprenoid Biosynthesis	<a href="#">HF610801</a>
ORF4	100	540 / 540	NODE_27_length_2659_cov_706.860840	1380..1919	Putative isopentenyl-diphosphate delta-isomerase	<a href="#">AFRH01000026</a>
aaiC	100	507 / 507	NODE_650_length_19027_cov_20.347137	16102..16608	Type VI secretion protein	<a href="#">cp003301</a>

# Issues to be aware of

- **Data quality**
- **Sequence bias**
  - Undetected/partly detected genes
- **Non-expressed genes**
  - Phenotype vs genotype problem
- **New variants and multiple variants of same gene**
- **Serotypes do not always correlate with phylogeny**
  - Combine with MLST
  - Enterobase is good for overview of serotypes and ST's
- **Consider coverage, ID, and positions**
- **Much output for VirulenceFinder**
  - Sort and select relevant genes

# Introduction to exercise

- ***Salmonella***
  - MLST
  - Serotype – Seqsero2
- ***E. coli***
  - Serotype – SerotypeFinder
  - Virulence - VirulenceFinder
- Note the results in sheets and discuss in group
- Any problems?

# Results - *Salmonella*

Strain ID	ST	Species	Subspecies	Serotype	O antigen	H1 antigen	H2 antigen	Antigen profile	Comments
SRR27241771	50	enterica	enterica	Saintpaul	4	e,h	1,2	4:e,h:1,2	No notes
SRR27241772	198	enterica	enterica	Kentucky	8	i	z6	8:i:z6	No notes
SRR27944993	3226	enterica	enterica	Paratyphi C or Choleraesuis or Typhisuis	7	c	1,5	7:c:1,5	<p><b>The predicted serotypes share the same general formula: 7:c:1,5 and can be differentiated by additional analysis</b></p> <p>Biochemical test necessary for differentiation – important to determine since Paratyphi C can cause invasive disease (enteric fever) and Choleraesuis is host adapted to pigs and Typhisuis is host restricted to pigs</p>
SRR27944994	86	enterica	enterica	Paratyphi B	4	b	1,2	4:b:1,2	<p><b>Detected the SNP in gene STM3356 that is associated with the d-tartrate nonfermenting phenotype characteristic of the typhoidal pathotype</b></p> <p>Paratyphi B (d-tartrate negative) can cause enteric fever and the Java variant (d-tartrate fermenting/positive) is less pathogenic and cause gastroenteritis</p>

# Results – *E. coli*

Strain ID	O	H	Serotype	eae genes	stx genes	stx subtype	Repeated genes	Comments
STEC-EQA-11-BB-10	O187	H28	O187:H28	none	stx2g-Out-S-8	stx2g	<b>gad, terC</b> , traT	traT – multiple variants in database and/or multiple copies of genes
STEC-EQA-11-BH-1	O80	H2	O80:H2	eae-e07-xi	stx2d-O55-5905 stx2d-OR-TS06-08	stx2d	afaB, cia, etsC, <b>gad</b> , iss, mchF, nleB, nleC, ompT, shiA, <b>terC</b> , traJ	<b>O80 is only detected in wzx – and with low coverage (812/1221) and ID% (96,43%).</b> Additional analysis needed to determine the O-type. afaB – bug (same reference sequence) cia, etsC, iss, mchF nleB, nleC, ompT, shiA - multiple variants in database and/or multiple copies of genes traJ – same gene with different reference sequences – clean up in database.
STEC-EQA-11-strain0017	O157	H7	O157:H7	eae-g01-gamma	stx1a-O157-FLY16, stx2c-O157-FLY16	stx1a, stx2c	astA, <b>gad</b> , nleB, tccP, <b>terC</b>	astA, nleB, tccP - multiple variants in database and/or multiple copies of genes

**gad and terC** – found in almost all *E. coli* - disregard



# Further reading and references

- ***Salmonella* serotyping and MLST**
  - White-Kauffmann-Le Minor serotype classification scheme [White-Kauffmann-Le Minor Scheme](https://www.pasteur.fr/sites/default/files/veng_0.pdf) ([https://www.pasteur.fr/sites/default/files/veng\\_0.pdf](https://www.pasteur.fr/sites/default/files/veng_0.pdf))
  - SeqSero2 – <http://denglab.info/SeqSero2>, <https://doi.org/10.1128/AEM.01746-19>
  - SeqSero1 – <http://www.denglab.info/SeqSero>, <https://doi.org/10.1128/jcm.00323-15>
  - Sequence bias SeqSero - <https://doi.org/10.1128/AEM.00614-20>
  - MLST and eBG - <https://doi.org/10.1371/journal.ppat.1002776>
  - Enterobase - <https://enterobase.warwick.ac.uk/>
  - SISTR - <https://doi.org/10.1371/journal.pone.0147101>

# Further reading and references

- ***E. coli* serotyping and virulence**
  - Pathotypes review - <https://doi.org/10.1093%2Ffemsre%2Ffuac031>
  - SerotypeFinder - <https://cge.food.dtu.dk/services/SerotypeFinder/>, <https://journals.asm.org/doi/10.1128/jcm.00008-15>
  - VirulenceFinder - <https://cge.food.dtu.dk/services/VirulenceFinder/>, <https://doi.org/10.1128/jcm.03617-13>, <https://doi.org/10.1128/jcm.01269-20>

# Acknowledgements

The creation of this training material was commissioned by ECDC to Statens Serum Institut (SSI) and The Danish Technical University (DTU) with the direct involvement of Egle Kudirkiene (SSI), Sofie Holtsmark Nielsen (SSI), Katrine Grimstrup Joensen (SSI) and Pernille Gymoese (SSI).